

(12)特許協力条約に基づいて公開された国際出願

02 FEB 2005

(19) 世界知的所有権機関
国際事務局(43) 国際公開日
2004 年 8 月 12 日 (12.08.2004)

PCT

(10) 国際公開番号
WO 2004/068399 A1

- (51) 国際特許分類: G06N 5/00
(21) 国際出願番号: PCT/JP2004/000971
(22) 国際出願日: 2004 年 1 月 30 日 (30.01.2004)
(25) 国際出願の言語: 日本語
(26) 国際公開の言語: 日本語
(30) 優先権データ:
特願2003-023342 2003 年 1 月 31 日 (31.01.2003) JP
(71) 出願人 (米国を除く全ての指定国について): 松下電
器産業株式会社 (MATSUSHITA ELECTRIC INDUS-
TRIAL CO., LTD) [JP/JP]; 〒5718501 大阪府門真市大
字門真 1 0 0 6 番地 Osaka (JP).
(72) 発明者; および
(73) 発明者/出願人 (米国についてのみ): 森川 幸治

(MORIKAWA, Koji) [/]. 大森 隆司 (OMORI, Takashi)
[/]. 大東 優 (OHIGASHI, Yu) [/]. 岡 夏樹 (OKA,
Natsuki) [/].

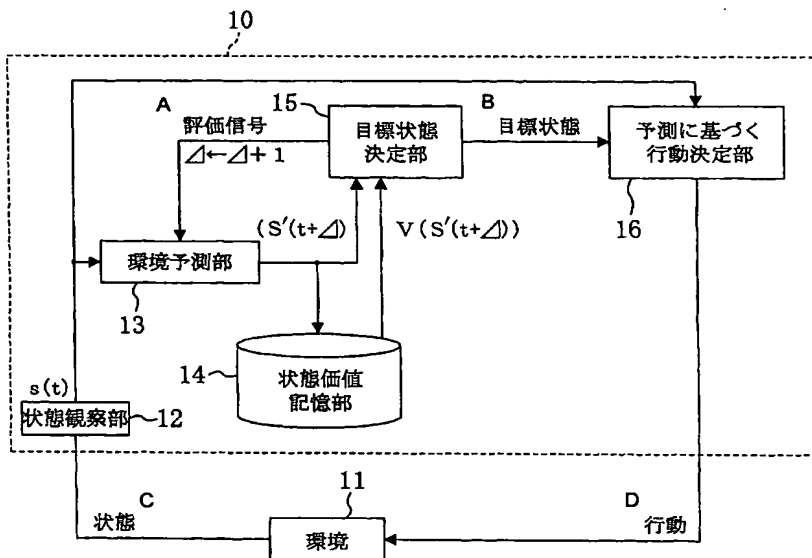
(74) 代理人: 前田 弘, 外 (MAEDA, Hiroshi et al.); 〒
5500004 大阪府大阪市靱本町 1 丁目 4 番 8 号 本町
中島ビル Osaka (JP).

(81) 指定国 (表示のない限り、全ての種類の国内保護が
可能): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR,
BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM,
DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU,
ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS,
LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA,
NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE,
SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US,
UZ, VC, VN, YU, ZA, ZM, ZW.

[続葉有]

(54) Title: PREDICTIVE ACTION DECISION DEVICE AND ACTION DECISION METHOD

(54) 発明の名称: 予測型行動決定装置および行動決定方法



- 12...STATE OBSERVATION SECTION
13...ENVIRONMENT PREDICTION SECTION
14...STATE VALUE STORAGE SECTION
15...TARGET STATE DECISION SECTION
A...EVALUATION SIGNAL
B...TARGET STATE
16...ACTION DECISION SECTION BASED ON PREDICTION
C...STATE
D...ACTION
11...ENVIRONEMNT

(57) Abstract: A predictive action decision device (10) includes a state observation section (12) which observes the state of an environment (11) and acquires a state value $s(t)$. An environment prediction section (13) predicts a state change of the environment (11) in future according to the state value $s(t)$. A target state decision section (15) references a state value storage section (14) and decides a future state appropriate for action decision as a target state. An action decision section (16) based on the prediction decides an action according to the target state decided.

(57) 要約: 予測型行動決定装置 (10) において、状態観察部 (12) は環境 (11) の状態を観察し、状態値 $s(t)$ を取得する。環境予測部 (13) は状態値 $s(t)$ に基づいて、環境 (11) の将来の状態変化を予測する。目標状態決定部 (15) は状態価値記憶部 (14) を参照して、行動決定のために適した将来状態を目標状態として

決定する。予測に基づく行動決定部 (16) は決定された目標状態を基にして、行動を決定する。



(84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), ヨーロッパ (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

添付公開書類:

- 国際調査報告書
- 補正書

2文字コード及び他の略語については、定期発行される各PCTガゼットの巻頭に掲載されている「コードと略語のガイダンスノート」を参照。

明 細 書

予測型行動決定装置および行動決定方法

技術分野

本発明は、外部から入力を受け、現在の状態から将来どのような状態に遷移するかを予測しながら、外部への出力を決定する行動決定に関する技術に属する。

背景技術

近年、産業上使用されるシステムは、年々複雑化しており、入力と出力との関係を予めプログラム等によって記述しておくことが困難になりつつある。このため、入力信号を処理して正しい出力を求める方法が必要になっており、このような入力から出力を決定する装置のことを、本願明細書中で「行動決定装置」と呼ぶ。また、特に、入力信号から将来の状態遷移を予測した上で、出力を求めるものを「予測型行動決定装置」と呼ぶ。

行動決定のための従来技術は、1) 現在の状態のみによって行動決定を行うもの、2) 過去の状態の遷移によって行動決定を行うもの、3) 将来の状態を予測して行動決定を行うもの、に分類される。

現在の状態のみによって行動決定を行う従来技術として、I F - T H E N ルールを用いるもの、ニューラル・ネットワークを用いるもの、テーブル参照方式などがある。これらは、現在の状態に対する行動が予め記述されており、入力から現在の状態を判別し、所定の記述を参照して行動を決定する。

ところが、現在の状態だけで正しい行動が決定できるとは必ずしも限らない。例えば対話型ロボットシステムにおいて、「いいですか」と聞かれたとき、それだけでは意味不明であり、それまでの状態遷移を参照することによってはじめて言葉の意味を理解することができる。すなわち、行動決定のために、過去の状態が必要となる場合もある。

また、現在や過去の状態だけでなく、将来の状態も考慮すべき場合もある。例

例えば、移動ロボットが障害物を回避するような場合には、まだ衝突していない段階では問題は発生していなくて、将来このままの移動方向や速度では衝突するという将来の状態遷移を考慮することによって、はじめて衝突前に回避の行動をとることができる。

将来の状態を考慮した行動決定に関する従来技術として、特許文献 1，2 に開示された技術がある。特許文献 1 では、環境から視覚センサや関節角度センサによって得た入力信号から、現在の画像データや関節角度データを状態として受ける。システムは、対象物に対する行動とその行動結果としての画像の変化をリカレント・ニューラル・ネットワークによって記憶し、同じような状況を受けた場合に、記憶している行動を再現する。この技術は、例えばロボットの自律的行動決定に応用されている。また特許文献 2 では、強化学習における行動決定技術が示されており、ある状態における価値と 1 ステップ前の状態における価値から、誤差を予測してその情報を行動決定に用いている。

また、特許文献 3 では、車両の安全走行制御のために、走行経路を推定し、現在の速度で走行を続けた場合に危険な地点がある場合に、その地点に到達する以前に運転者の技量に応じて安全な速度に車両速度を制御する技術が開示されている。

(特許文献 1) 特開 2002-59384 号公報

(特許文献 2) 特開 2002-189502 号公報

(特許文献 3) 特開平 7-306998 号公報

解決課題

上述のように、特許文献 1 では、現在状態から自己の行動によってどのような状態に変化するかをリカレント・ニューラル・ネットワークを用いて予測している。そして、その予測結果に応じて、状態と対で記憶された行動が決定される。

しかしながら、特許文献 1 では、自己の行動に起因する過去の状態変化をリカレント・ニューラル・ネットワークを用いて学習しているに過ぎず、自己の行動と関係のない環境の変化については、予測も考慮も何らなされていない。また、ある時点での行動決定は、現在の状態とその 1 ステップ後の状態予測に基づいて

行われているが、この1ステップ後の状態が行動決定のために重要であるとは必ずしもいえず、したがって、この将来の状態予測が行動決定にとって適切なものとはいえない。

また、特許文献2に関しても、現在の状態と1ステップ後の状態の予測価値のみから決定された行動が、必ずしも望ましいものとは限らない。例えば、移動ロボットが走ってくる車をよけたいときに、移動ロボットの速度が車に対してかなり遅い場合は、何ステップも前から回避行動を起こさないと車とぶつかってしまうように、将来を見越して行動を決定すべき場合には、1ステップ後だけでなくさらに将来の状態も考慮して行動を決定する必要がある。また、例えば上述の移動ロボットが走ってくる車をよける場合のように、現在や1ステップ後の価値を見ても変化はみられないが、何ステップも先に決定的な自体が発生するときは、現在や1ステップ後の価値に基づく行動は、無駄な行動につながる可能性もある。

また、特許文献3では、地図情報と自車位置情報とを用いて、車が将来進む経路を推定した上で、現在の速度で推定経路を走行した場合に危険な地点が存在するとき、各地点の目標車速を設定した車速計画を作成する。すなわち、車は道路を走行する、という前提の元に、その道路情報を用いて、対処すべき目標となる危険な地点の有無を判断している。ところが、道路情報のような将来状態の情報が予め与えられていない場合には、行動決定を行うための適切な目標を定めること自体が、決して容易ではない。また、この技術では、様々な未経験の事態に対して行動を決定することが、極めて困難である。

発明の開示

本発明は、予測型行動決定装置において、行動決定のために将来の状態予測をより適切に行い、行動決定の精度や能力を向上させることを目的とする。

本発明は、次の点に着目している。すなわち、環境の状態の中には、自己（予測型行動決定装置）の行動（＝出力）に関連するものと、自己の行動に関係なく変化するものとがあり、前者については、長い将来にわたっての状態予測は困難であるが、その一方で、後者については、1ステップ先だけではなく長い将来に

わたっての予測も比較的容易に行うことができる。

すなわち、本発明では、環境を観察して、環境の将来の状態変化を予測し、予測した各将来状態に係る状態価値を参照しつつ、行動決定のために適した将来の状態を、目標状態として決定する。そして、この決定した目標状態を基にして、自己の行動を決定する。これにより、将来状態に係る状態価値を参照して決定された目標状態を基にして、行動が決定されるので、従来よりも行動決定の精度が向上する。

図面の簡単な説明

図 1 は本発明、および各実施形態を説明するための課題を示す図である。

図 2 は政策に基づく行動決定に用いられる状態価値と行動基準の例である。

図 3 は本発明に係る、予測に基づく行動決定を示す図である。

図 4 は本発明の第 1 の実施形態に係る予測型行動決定装置の構成を示すブロック図である。

図 5 は図 4 の予測型行動決定装置の動作を示すフローチャートである。

図 6 は本発明の第 2 の実施形態に係る予測型行動決定装置の構成を示すブロック図である。

図 7 は図 6 の構成における価値変換部の動作を示すフローチャートである。

図 8 は予測ステップ数に伴う状態価値の変化の一例を示すグラフである。

図 9 は本発明の第 3 の実施形態における予測に基づく行動決定部の内部構成を示す図である。

図 10 は図 9 の構成の動作を示すフローチャートである。

図 11 は本発明の第 4 の実施形態における環境予測部の内部構成を示す図である。

図 12 は図 11 の構成の動作を示すフローチャートである。

図 13 は本発明の実施形態に係る学習方法と従来の強化学習のシミュレーション結果を示すグラフである。

図 14 は本発明の実施形態に係る学習方法と従来の強化学習のシミュレーション結果を示すグラフである。

図 1 5 は本発明の実施形態に係る学習方法と従来の強化学習のシミュレーション結果を示すグラフである。

図 1 6 は状態価値の学習過程を概念的に示す図である。

図 1 7 は学習回数に応じた状態価値曲線の変化を概念的に示す図である。

図 1 8 は本発明の第 2 の実施形態の変形例に係る予測型行動決定装置の動作を示すフローチャートである。

図 1 9 は本発明の第 2 の実施形態の変形例における行動選択基準の一例を概念的に示す図である。

発明を実施するための最良の形態

本発明の第 1 態様によれば、所定の環境について状態を観察し、状態値を取得する状態観察部と、前記環境の各状態に係る状態価値を記憶する状態価値記憶部と、前記状態観察部によって取得された状態値に基づいて、前記環境の将来の状態変化を予測する環境予測部と、前記環境予測部によって予測された各将来状態の中から、前記状態価値記憶部に格納された前記各将来状態の状態価値に基づいて、行動決定のために適した将来の状態を、目標状態として決定する目標状態決定部と、前記目標状態決定部によって決定された目標状態を基にして、自己の行動を決定する第 1 の行動決定部とを備えた予測型行動決定装置を提供する。

本発明の第 2 態様によれば、前記環境予測部は、前記環境の将来の自己の行動に影響されない状態変化を予測する第 1 態様の予測型行動決定装置を提供する。

本発明の第 3 態様によれば、前記目標状態決定部は、状態価値が極大となる将来状態を目標状態として決定する第 3 態様の予測型行動決定装置を提供する。

本発明の第 4 態様によれば、前記状態価値記憶部に記憶された状態価値を学習によって更新する状態価値更新部を備え、前記目標状態決定部は、前記各将来状態の中で、状態価値が前記状態価値更新部によってすでに更新されているものを、前記目標状態として決定する第 1 態様の予測型行動決定装置を提供する。

本発明の第 5 態様によれば、前記目標状態決定部は、前記状態価値記憶部から得た状態価値を現在からのステップ数に応じて割り引いて用いる第 1 態様の予測型行動決定装置を提供する。

本発明の第 6 態様によれば、前記状態価値記憶部は、自己を含む状態に係る状態価値を記憶するものであり、当該予測型行動決定装置は、前記環境予測部によって予測された、自己を含まない将来状態について、前記状態価値記憶部に記憶された状態価値を基にしてその状態価値を求め、前記目標状態決定部に与える価値変換部を備えた第 1 態様の予測型行動決定装置を提供する。

本発明の第 7 態様によれば、所定の行動基準に基づいて自己の行動を決定する第 2 の行動決定部と、前記第 1 および第 2 の行動決定部によって決定された行動を第 1 および第 2 の行動候補として受け、これら第 1 および第 2 の行動候補のうちのいずれか一方を実際の行動として選択する行動選択部とを備えた第 1 態様の予測型行動決定装置を提供する。

本発明の第 8 態様によれば、前記目標状態決定部は、目標状態を決定できたか否かを示す選択信号を前記行動選択部に与えるものであり、前記行動選択部は、前記選択信号が、目標状態を決定できたことを示すときは、前記第 1 の行動候補を選択する一方、目標状態を決定できなかったことを示すときは、前記第 2 の行動候補を選択する第 7 態様の予測型行動決定装置を提供する。

本発明の第 9 態様によれば、前記第 1 の行動決定部は、前記状態値を受け、この状態値が表す現在状態からその前ステップにおける状態と行動を検出する行動付状態変化検出部と、前記行動付状態変化検出部によって検出された、現在状態並びにその前ステップにおける状態および行動の組合せを、状態変化として記憶する行動付状態変化記憶部と、前記行動付状態変化記憶部から、現在状態から目標状態までの期間の状態変化の履歴を検索し、この検索結果を基にして、行動を決定する行動計画部とを備えた第 1 態様の予測型行動決定装置を提供する。

本発明の第 10 態様によれば、前記行動計画部は、前記状態変化記憶部の検索の際に、目標状態から現在の状態に向かって後ろ向き探索を行う第 9 態様の予測型行動決定装置を提供する。

本発明の第 11 態様によれば、前記環境予測部は、前記状態値を受け、この状態値が表す現在状態からその前ステップにおける状態を検出する状態変化検出部と、前記状態変化検出部によって検出された現在状態およびその前ステップにおける状態の組合せを状態変化として記憶する状態変化記憶部と、前記状態変化記

憶部から現在状態の後の状態を予測する状態予測部とを備えた第1態様の予測型行動決定装置を提供する。

本発明の第12態様によれば、予測型行動決定装置において、自己の行動を決定する方法として、所定の環境について状態を観察して、状態値を取得する第1のステップと、取得した状態値に基づいて、前記環境の将来の状態変化を予測する第2のステップと、予測した各将来状態の中から、前記各将来状態の状態価値を参照しつつ、行動決定のために適した将来状態を目標状態として決定する第3のステップと、決定した目標状態を基にして、自己の行動を決定する第4のステップとを備えた行動決定方法を提供する。

本発明の第13態様によれば、予測する状態変化は、前記環境の将来の自己の行動に影響されない状態変化である第12態様の行動決定方法を提供する。

本発明の第14態様によれば、前記第3のステップにおいて、前記各将来状態の中で、状態価値が極大となるものを、前記目標状態として決定する第12態様の行動決定方法を提供する。

本発明の第15態様によれば、前記予測型行動決定装置は、前記環境の各状態に係る状態価値を学習によって更新するものであり、前記第3のステップにおいて、前記各将来状態の中で状態価値がすでに更新されているものを、前記目標状態として決定する第12態様の行動決定方法を提供する。

まず、本発明に関する基本的な概念について説明する。

図1は課題の例を示す図である。図1では、座標(0, 0) - (1, 1)の空間1において、ボールBが直進運動をしている。空間1の上下左右の壁に当たったとき、ボールBは反射する。パドルPは左右方向にのみ動くことができるものとする。ボールBの状態は、位置座標(B_x , B_y)と進む方向 B_t によって表現され、パドルPの状態は、位置座標(P_x , P_y)によって表現される。ただし、 $P_y = 0$ で固定されている。

各時刻ステップにおいて、パドルPの操作MPとして、{LEFT (左に動く)、RIGHT (右に動く)、STAY (そのまま)}のうちいずれか1つを選択する。そして、ボールBをパドルPによって受けることができたとき、正の報酬が得られるものとし、一方、ボールBをパドルPによって受けることができな

かったとき、負の報酬が得られるものとする。そしてここでのタスクは、ボールBをパドルPによって受ける回数をできるだけ増やし、累積で得られる報酬をより多くすることである。

これを行動決定問題として見た場合、各ステップごとに、入力としてボールBおよびパドルPの状態 (B_x , B_y , B_t , P_x) が与えられたとき、より多くの報酬が得られるように、パドルPの操作MPを選択する、ということになる。

このような課題に対して、例えば、状態価値および行動基準を用いて、行動を決定する方法がすでに知られている。このような方法を、本願明細書では「政策 (Policy) に基づく行動決定」と呼ぶ。

図2は政策に基づく行動決定に用いられる状態価値と行動基準の例を示す図である。ここでの「状態価値」は、外部の状態を評価して得た価値のことである。図2の例では、図1の空間1が 8×8 のセルに分割されており、各セルには、ボールBがその位置にあるときの状態価値が示されている。例えばボールBがセルCLの位置にあるとき、状態価値は「3」である。図2の例では、ボールBが空間1内の下面に到達したときに報酬がもらえるか否かが決定されることから、下面に近い位置ほど状態価値が高くなっている。状態価値は、事前に与えられたり、学習によって獲得されたりする。

また、各セルには、ボールBがその位置に来たときどの行動を取るべきか、という行動基準が記述されている。図2の例では、図1と対応させて考えると、パドルPがボールBの下あたりにあり、ボールBが左下に向かって落ちてきているときに、ボールBの位置に対応したセルCLにおける行動基準として、左に動く行動に対して0.5、動かない行動に対して0.4、右に動く行動に対して0.1、という値が割り当てられている。この値を基にして、例えば最大値をもつ行動を選択したり、行動選択確率として計算したりして、最終的な行動が決定される。例えば最大値をとる行動を選択する場合には、セルCLでは、最大値0.5を持つ「左に動く」が選択される。

このように、政策に基づく行動決定では、現在のボールBの状態 (位置) に対応する状態価値が参照され、その状態に対する行動基準によって行動が決定される。すなわち、行動が現在の状態を考慮して決定されるので、例えばボールBの

速度が早い場合など、事前に行動を決定する必要があるような場合等には対応できなかった。

これに対して、本発明では、「予測に基づく行動決定」を行う。すなわち図 3 (a) に示すように、現在のボール位置 B 1 に対して、複数のステップにわたる将来のボール位置を予測し、行動決定に適した目標状態となるボール位置 B 2 を決定し、この目標状態を基にして行動決定を行う。

図 3 (b) は予測された各ボール位置に対する状態価値が格納されたテーブル、図 3 (c) は予測ステップ数に対する状態価値の変化を示すグラフである。図 3 (c) から分かるように、予測ステップ数を増やすにつれて、ボール B は下面に向かって移動することになり、これとともに状態価値は徐々に増加する。そして、下面に到達してからはボール B は上方に向かうため、状態価値は減少に転じる。このため、状態価値は、予測ステップ数 5 のとき、最大値「8」になる。ここで、例えば最大の状態価値を取った予測ステップ数と予測状態を行動決定に用いるものとする、位置 B 2 が目標状態となる。これは、ボール B の状態でいうと、一番下面に近いときが目標状態として設定されることに相当する。このように、将来の、自己の行動に影響されない状態変化を予測して、予測した状態の中で行動決定のために適したものを目標状態として設定して行動決定を行うことによって、現在の行動をよりの確に決定することができる。

以下、本発明の実施の形態について、図面を参照して説明する。なお、以下の説明では、図 1 の課題をタスクとして扱うものとする。図 1 の課題は、ピンポンのようにボールが移動してそれを跳ね返すと報酬が得られるものであるが、これは、サッカーゲームにおいてボールを取りに行く行動、移動ロボットが接近物を回避する行動等と同様の設定とみなすことができる。接近物の回避では、接近物の動きを予測してこれを回避することを目標状態として設定することによって、行動を決定できる。この他にも、ロボットが物を受け取る際の行動決定なども、同様の課題となる。

(第 1 の実施形態)

図 4 は本発明の第 1 の実施形態に係る予測型行動決定装置 10 の構成を示す。この予測型行動決定装置 10 は、所定の環境 11 について状態を観察し、環境 1

1 の、自己の行動に影響されない状態変化を予測した上で、目標状態を定め、環境 1 1 に対する行動を決定する。

状態観察部 1 2 は環境 1 1 の状態を観察し、現在の状態を表す状態値を取得する。ここでは、状態値を $s(t)$ と表す。図 1 の課題については、ボール B およびパドル P の状態 (B_x , B_y , B_t , P_x) が状態値 $s(t)$ として得られる。

環境予測部 1 3 は状態観察部 1 2 によって取得された状態値 $s(t)$ に基づいて、環境 1 1 の将来の状態変化を予測する。ここでは、ボール B の座標 (B_x , B_y) を状態 s' として予測するものとする。すなわち、ボール B は基本的には直進し、壁で反射するので、将来の座標は、現在のボール B の座標と角度から解析的な計算によって予測することができる。またこの予測には、パドル P の操作 MP を考慮する必要はない。すなわち、当該装置 1 0 の行動に影響されない、将来の状態変化を予測することになる。本実施形態では、環境予測部 1 3 はこのような予測機能を予め備えているものとする。なお、このような予測機能は学習等によって得ることも可能であり、これについては後述する。

状態価値記憶部 1 4 は図 3 (b) に示すように、各状態 s' すなわちボールの位置 (B_x , B_y) について状態価値をそれぞれ格納している。ここでは、状態 s' に対する状態価値を $V(s')$ と記述する。本課題では、ボール B が下の壁に当たるとき、報酬がもらえるか否かが判定されるので、下の壁に近い位置ほど、より高い状態価値が設定されている。なお、状態価値は全ての状態について設定されているのが好ましいが、これは実際には困難なので、例えば特定の部分状態にのみ状態価値を設定しておいてもよい。

また、状態 s' の要素にボール B の進む方向 B_t を含めてもよく、この場合は、状態価値記憶部 1 4 に、(B_x , B_y , B_t) について状態価値をそれぞれ格納しておけばよい。

目標状態決定部 1 5 は、環境予測部 1 3 による予測結果を基にして、行動の決定のために適した将来の状態を目標状態として決定する。ここでは、環境予測部 1 3 によって複数ステップにわたって予測された将来の状態 $s'(t + \Delta)$ について、それぞれ、状態価値記憶部 1 4 を参照して状態価値 $V(s'(t + \Delta))$

を求め、求めた状態価値から目標状態を決定する。

第1の行動決定部としての予測に基づく行動決定部16は、目標状態決定部15によって決定された目標状態に対して、現在どのような行動を取るべきかを決定する。ここでは、ボールBの座標とパドルPの座標が一致したとき報酬がもらえるので、目標状態から、パドルPがどの位置に向かえば報酬をもらえる可能性が高くなるを識別し、その位置に向かうように、操作MPを決定する。

以下、図5のフローチャートを参照して、図4の予測型行動決定装置10の動作を説明する。まず、状態観察部12によって、環境11から、現在の状態を表す状態値 $s(t)$ が取得される(S11)。そして、何ステップ先の予測を行うかを変数 Δ で定義し、その初期値として1を与える(S12, S13)。なお、ステップS17で所定の条件を満たすまで、または、 Δ が所定値 n を超えるまで(S19)、 Δ をインクリメントしながら、以下のステップS14～S16を繰り返し実行する。

ステップS14において、環境予測部13は、現在状態から Δ ステップ先の状態 $s'(t+\Delta)$ を予測する。例えば、ボールBが壁に当たらないときは、 B_x , B_y が1ステップ分変化し、方向 B_t は変化しない。そしてステップS15において、状態価値記憶部14から、ステップS14で予測された状態 $s'(t+\Delta)$ の状態価値 $V(s'(t+\Delta))$ が抽出される。

ステップS16において、目標状態決定部15は、状態価値記憶部14から出力された状態価値 $V(s'(t+\Delta))$ を評価する。ここでは、評価条件として、所定の閾値を超えているか否かを判定する。そして、条件を満たす、すなわち状態価値が所定の閾値を越えていると判定したときは(S17でYES)、ステップS18に進み、状態 $s'(t+\Delta)$ を目標状態として、予測に基づく行動決定部16に与える。

一方、条件を満たさないときは(S17でNO)ステップS19に進み、予測ステップ数の判定を行う。すなわち、 Δ が所定値 n を下回るときは(S19でNO)ステップS13に戻り、 Δ をインクリメントした後、同様の処理を行う。一方、 Δ が所定値以上のときは(S19でYES)ステップS1Aに進み、目標状態が決定できなかった旨を予測に基づく行動決定部16に通知する。

予測に基づく行動決定部 16 は、目標状態決定部 15 からの出力を受けて、次の行動を決定する。しかし、目標状態が決定できなかった通知を受けた場合は、例えばランダムに、行動を決定する。

このように本実施形態によると、環境 11 の状態変化の予測結果から、状態価値を参照して、行動決定のために適した将来の状態が目標状態として決定される。そして、この目標状態を基にして、自己の行動が決定される。このため、従来の将来を予測した行動決定よりも、行動決定の精度が格段に向上する。また、従来技術のように状態と行動との関係を予め記述していなくても、行動が決定できるので、簡易な構成により、様々な未経験の事態に対しても行動決定可能になる。また、環境 11 の将来の、当該装置 10 の行動に影響されない状態変化を予測しているので、1 ステップ先だけではなく長い将来にわたっての予測も比較的容易に、精度良く、行うことができる。

なお、本実施形態では、状態価値が所定の閾値を越えたときの状態を目標状態として決定するものとしたが、目標状態の決定は、これ以外にも、様々な方法が考えられる。例えば、予測したステップ内で状態価値が最大となった状態や、状態価値が極大になった状態を、目標状態として決定してもよい。また、前ステップとの状態価値の差分が所定値よりも小さくなったときの状態を目標状態として決定してもよい。あるいは、現在からのステップ数に応じて、状態価値を割り引いて評価するようにしてもよい。

なお、本実施形態では、状態価値がテーブル形式で記憶されているものとしたが、この代わりに、ニューラル・ネットワーク等の関数近似手法を適用してもよい。この場合、現在の状態が入力されたとき、状態価値または報酬の期待値が出力されるように学習がなされている必要がある。

(第 2 の実施形態)

図 6 は本発明の第 2 の実施形態に係る予測型行動決定装置 10 A の構成を示すブロック図である。図 6 において、図 4 と共通の構成要素には図 4 と同一の符号を付している。

まず、価値変換部 21 について説明する。価値変換部 21 は、状態価値記憶部 14 A が当該装置 10 A の行動によって変化する状態を含む状態に係る状態価値

を記憶しているとき、環境予測部 13 によって予測された、当該装置 10A の行動によって変化する状態を含まない将来状態 $s' (t + \Delta)$ の状態価値 $V (s' (t + \Delta))$ を、状態価値記憶部 14A に格納された状態価値を基にして求め、目標状態決定部 15A に与える。本実施形態では、状態価値記憶部 14A は、パドル P の位置 P_x も含めた状態 (B_x, B_y, B_t, P_x) についてそれぞれ状態価値を記憶しているものとし、環境予測部 13 からは将来状態 (B_x, B_y, B_t) が出力されるものとする。

図 7 のフローチャートを参照して、価値変換部 21 の動作を説明する。まず、環境予測部 13 によって将来状態 $s' (t + \Delta)$ が予測されたとき、状態価値記憶部 14A をこの状態 $s' (t + \Delta)$ によって検索する (S21)。そして、3 個の要素 (B_x, B_y, B_t) が一致する状態について、その状態価値の組 $V' (s' (t + \Delta))$ を状態価値記憶部 14A から抽出する (S22)。ここで、パドル P の座標 P_x が 8 通りであるとする、8 個の状態価値が探索結果として出力される。

そして、ステップ S22 において抽出された状態価値を比較し、その最大値を求め (S23)、得られた最大値を状態 $s' (t + \Delta)$ の状態価値 $V (s' (t + \Delta))$ として目標状態決定部 15A に出力する。

また、第 1 の実施形態では、状態価値記憶部 14 に格納された状態価値は事前に与えられるものとしたが、本実施形態では、強化学習と呼ばれる手法によって自動的に学習させるものとする。「強化学習」とは、参考文献 1 (サットン (R. S. Sutton), パート (A. Barto) 著, 「強化学習入門 (Reinforcement Learning: An Introduction)」, (米国), A Bradford Book, The MIT Press, 1998 年 3 月) で知られる学習方法の 1 種であり、環境からの報酬信号と試行錯誤によって学習を行う手法であり、報酬の最大化を目的として行動決定を行う学習アルゴリズムのことである。

本実施形態では、強化学習の中で、actor-critic 学習という手法を利用する (参考文献 1 の p. 151 ~ 153 参照)。actor は、どの状態でどの行動を取るべきかを記述する行動決定方策 policy を持ち、状態 s で、とりうる行動 $a_1 \sim a_n$ に対して選択確率が算出されて決定される。crit

$i c$ の V には状態価値と呼ばれる、ある状態でどれぐらい報酬がもらえそうかの期待値を示す値が格納される。状態価値は、報酬から計算される TD 誤差によって更新される。TD 誤差 δ は、

$$\delta = r(t+1) + \gamma (V(s(t+1)) - V(s(t)))$$

で計算される。ここで $r(t+1)$ は報酬、 γ は割引率 (discount rate) と呼ばれている。この TD 誤差によって、 $c r i t i c$ 内部の状態価値テーブル $V(s)$ と、 $a c t o r$ の持つ $P o l i c y$ の更新を行う。

$a c t o r - c r i t i c$ 法を利用して、状態価値記憶部 14 A に記憶される状態価値を更新することによって、状態価値が不明な状態であっても、予測によって行動を決定することができる。状態価値更新部 22 は $c r i t i c$ テーブル更新則によって、状態価値記憶部 14 A 内に記憶された状態価値を更新する。

ここで、状態価値の学習過程について、図 16 を用いて概念的に説明する。状態価値を学習によって獲得する場合、最初は状態価値は与えられておらず、上述した強化学習などの学習法に従って、徐々に状態価値が学習される。

図 16 では、図 2 と同様に、図 1 の空間 1 が 8×8 のセルに分割されており、各セルには、ボール B がその位置にあるときの状態価値が示されている。同図中、(a) は学習前の状態であり、どのセルも状態価値として初期値「0」が与えられている。この場合は、どの位置においても適切な行動は定義されておらず、価値に基づいて行動を決定することはできない。

図 16 (b) は学習途中の状態である。状態価値は報酬を受ける毎に更新されていくが、ここではパドル P がボール B を打ったときに正の報酬が得られるという設定なので、報酬はボール B が一番下にあるときにももらえることになる。そして、強化学習によると、状態価値は報酬の伝播のような計算によって算出されるため、図 16 (b) に示すように、報酬が得られる位置付近、すなわち下の方だけについて、状態価値が算出されている。この場合は、状態価値が学習されている位置にボール B が来たときはその状態価値に応じて行動が決定できるが、状態価値が依然として初期値のままの位置では、行動が決定できない。

そして、十分な学習期間が取られた後には、図 16 (c) のように、全ての位置に対して状態価値が算出され、どの位置においても状態価値に基づく行動が可

能になる。

また、状態価値の学習過程は、図 17 のようにも表すことができる。図 17 では、図 3 (c) と同様に、予測ステップ数に対する状態価値の変化をグラフに表している。ボール B の将来の位置を予測できる場合は、図 16 のような状態価値を参照することによって、横軸に予測ステップ数、縦軸に状態価値をとったときの状態価値の曲線を描くことができる。

上述したように、学習が進むに従って、状態価値はその値を増やし、また報酬がもらえるタイミング g から何ステップも前の状態に伝播していく。すなわち、学習回数が少ないときは、報酬がもらえるタイミング g の直前の状態においてのみ状態価値が学習され、学習が進むにつれて、タイミング g よりもだいぶ前の状態においても状態価値が与えられる。すなわち、状態価値曲線は、学習回数が増加するにつれて、曲線 c_1 、曲線 c_2 、曲線 c_3 のように次第に変化していく。

ここで例えば、曲線 c_3 のように状態価値が学習されたとすると、領域 e_1 は状態価値が割り当てられており、行動が状態価値によって決定できる、いわゆる既学習領域となる。一方、領域 e_2 は、まだ状態価値が学習されておらず、したがって状態価値による行動決定ができない、いわゆる未学習領域となる。曲線 c_3 が図 16 (b) の状態に対応しているものと考え、未学習領域 e_2 はボール B がまだ上方にある場合に該当する。

学習が十分に進んだ場合は、未学習領域 e_2 はなくなり、既学習領域 e_1 のみになる。これは、例えば図 16 (c) のような状態に対応していると考えられる。

このように、学習が進むと、各状態に状態価値が割り当てられていき、状態学習が与えられた既学習領域 e_1 が次第に大きくなるとともに、未学習領域 e_2 が次第に小さくなる。これは、強化学習の特徴の 1 つともいえる点である。このため、将来状態を複数ステップ予測したときに、状態価値が与えられている既学習領域 e_1 に達するまでにかかるステップ数が、学習が進むにつれて、徐々に少なくなっていくことになる。

図 8 は強化学習によって学習された状態価値記憶部 14 A の状態価値に基づき、価値変換部 21 が求めた状態価値の変化を示すグラフである。図 8 において、

縦軸は算出された状態価値、横軸は予測のステップ数である。また図 8 において、問題設定は図 1 と同様であるが、状態空間は図 2 よりも細かく分割されているため、予測ステップ数も多くなっている。図 8 のグラフでは、予測ステップ数の増加に伴って状態価値が増加しており、80 ステップあたりでピークを迎えている。一般に、状態価値は、図 8 のように、未学習領域から既学習領域に入ると増加し始め、その後に減少に転じる。この減少に転じるタイミングが、現在状態から予測される最も意義のある状態、すなわち目標状態に相当すると判断できる。

目標状態設定部 15 A は、第 1 の実施形態と同様に動作する。ただし、目標状態を決定できたか否かを示す選択信号を、行動選択部 24 に与える点が異なっている。予測に基づく行動決定部 16 A は、第 1 の実施形態と同様に動作して行動を決定し、決定した行動を第 1 の行動候補として行動選択部 24 に出力する。

第 2 の行動決定部としての政策に基づく行動決定部 23 は、所定の行動基準としての、強化学習における actor-critic 手法によって学習された政策に基づいて行動を決定し、決定した行動を第 2 の行動候補として行動選択部 24 に出力する。

行動選択部 24 は、予測に基づく行動決定部 16 A から受けた第 1 の行動候補と、政策に基づく行動決定部 23 から受けた第 2 の行動候補とのうち、いずれか一方を、環境 11 への実際の行動として選択する。この選択には、目標状態決定部 15 A から受けた選択信号を利用する。すなわち、選択信号が目標状態を決定できなかったことを示すときは、予測に基づく行動決定部 16 A から受けた第 1 の行動候補は意味がないと考えて、政策に基づく行動決定部 23 から受けた第 2 の行動候補を実際の行動として選択する。そうでない場合は、予測に基づく行動決定部 16 A から受けた第 1 の行動候補を実際の行動として選択する。

このように本実施形態によると、状態価値を強化学習によって求めることによって、あらかじめ与えることが困難であった状態価値が自律的に獲得され、予測型行動決定が容易に実現できる。また、学習装置としての観点から本装置を見た場合は、報酬が与えられる時点を予測することになり、初期の未学習の状態が多い場合にも行動の指針を得ることができ、学習の効率が向上する。

以下、従来の強化学習 (Reinforcement Learning ; 以下、「RL」と記載)と、

本実施形態に係る学習方法 (Prediction-based Reinforcement Learning; 以下「PRL」と記載) のシミュレーションによる比較結果をいくつか示す。このシミュレーションでは、ボールが打てたときの報酬は1.0、打てなかったときの報酬は-1.0、またパドルが右または左に動いたときの報酬は-0.01とした。また、PRLでは、最初に状態価値記憶部14A内の状態価値をある程度作成するために3000試行 (ボールが下面に当たるまでを1試行といい、以下epochと記載) は、RLと同様の処理で行動決定を行った。各epochにおけるボールの初期位置はランダムに設定した。

図13はRLとPRLそれぞれの累積報酬の変化を表すグラフである。横軸はepoch数、縦軸は累積報酬である。最初の3000epochでは、両手法の累積報酬にはほとんど差はない。これは、この期間では、PRLはRLと同じ行動決定方法と学習法で動いているからである。学習の3000epoch以降では、PRLの方が、環境予測部13と予測に基づく行動決定部16Aとの利用によって、RLよりも良い性能を示している。これは、PRLでは、まだ学習が収束していない (すなわち、状態価値記憶部14A内の状態価値が完全には生成されていない) ときでも、予測によって適切な行動が決定できるので、ボールを打ち返す確率がRLよりも高くなるからである。また、PRLはボールを打ち返す位置の予測に基づき行動を決定しているため、不必要な行動をする割合が低くなっていることもその原因の一つと考えられる。

図14は学習100epoch毎に評価を行った結果を示すグラフである。横軸はepoch数、縦軸は100epochごとに行われた性能評価の結果であり、ボールが打てた割合である。ここでのRLの行動決定規則には、ソフトマックス (softmax) と呼ばれる行動戦略を用いている。これは、各状態の各行動に付与された値に応じて確率的に行動決定を行う方法である。初期の3000epochの学習後、PRLのパフォーマンスが、RLと比べて大きく向上している。この結果は、まだボールを打ち返す付近の状態しか学習が進んでいない (すなわち、状態価値が決定されていない) RL学習の初期段階に、PRLが環境予測部13と予測に基づく行動決定部16Aとの利用によって、学習が進行していない状態に対しても適切な行動を決定できたことが大きな要因である。

図 15 は R L の行動決定規則を greedy 戦略に変更した場合のグラフである。図 14 と比べて、R L の性能は全体的に見ると向上している。しかし、学習していない状態に評価時に遭遇すると性能は急激に落ちる。R L は、P R L に比べると、精度が悪く安定した動作を示していないことが分かる。

なお、本実施形態では、状態価値をテーブルとして表現して強化学習により学習するものとしたが、ニューラル・ネットワークによって学習させることもできる。この場合、ニューラル・ネットワークの汎化能力によって未経験の状態に対する価値を出力することが期待される。しかし、この手法は、状態価値に不連続性が少ないときに有効であると考えられる。

(第 2 の実施形態の変形例)

また、学習との組合せを考えた場合、状態価値が未学習の状態から既学習の状態に切り換わる場所を目標状態として決定する方法も考えられる。この方法によると、学習がさほど進んでいない状況において、既学習領域では状態価値に基づいて行動決定する一方、未学習領域では、既学習領域に含まれた将来状態を目標状態として決定して、行動決定を行うことができる。

上述の第 2 の実施形態では、図 8 に示すグラフにおいて、状態価値が減少に転じるタイミングを目標状態として決定した。ところが、このようにして目標状態を決定できるためには、状態価値が十分な量になるまで蓄積される必要があり、問題によっては、それまでに膨大な学習回数が必要になる。

これに対して、本変形例では、未学習領域から既学習領域に入るタイミングを、目標状態として決定する。これにより、未学習領域が多い学習初期でも、未学習領域では目標状態を基にして行動を決定し、既学習領域では学習によって行動決定することができる。図 17 でいうと、未学習領域 e 2 と既学習領域 e 1 とを区別して、行動決定を行うことになる。

本変形例は、図 6 の予測型行動決定装置 10A の構成において、目標状態決定部 15A および行動選択部 24 の動作を変更することによって、実現される。以下、図 18 のフローチャートを用いて、その動作を説明する。

まずステップ S 51 において、目標状態決定部 15A は、状態観察部 12 によって観察された現在の状態 S (例えばボールの現在位置) から、現在の状態価値

Aを、状態価値記憶部14Aおよび価値変換部21を介して算出する。次にステップS52において、ステップS51で算出された状態価値Aは、状態価値更新部22によって既に更新されたことがあるか否か、すなわち既学習か否かを判断する。既学習であるときはステップS53に進み、そうでないときはステップS54に進む。

ステップS53では、現在の状態はすでに学習が進み状態価値が与えられているので、目標状態決定部15Aは行動選択部24に対して、政策に基づく行動決定部23が決定した第2の行動候補を選択するよう、選択信号を送る。これにより、政策に基づく行動決定部23によって決定された行動が、環境11への実際の行動として選択される。

一方、ステップS54では、現在の状態はまだ状態価値が未学習であると考えられるので、環境予測部13、価値変換部21および目標状態決定部15Aのループ動作によって、将来状態が予測されて、目標状態が決定される。ここでは、予測した将来状態が、既に状態価値が学習されている既学習領域に達したとき、その将来状態を目標状態として決定する。決定された目標状態は、予測に基づく行動決定部16Aに送られる。

そしてステップS55において、目標状態決定部15Aは行動選択部24に対して、予測に基づく行動決定部16Aが決定した第1の行動候補を選択するよう、選択信号を送る。これにより、予測に基づく行動決定部16Aによって決定された行動が、環境11への実際の行動として選択される。

このように、目標状態決定部15Aにおける目標状態の決定基準と、行動選択部24における行動選択基準とを変更することによって、学習が進んでいない段階であっても、政策に基づく行動決定を行うことができない未学習領域において、予測に基づく行動決定を行うことによって、より広い領域で、行動決定を行うことが可能になる。図19は本変形例における行動選択を概念的に示した図である。図19に示すとおり、学習回数が十分でないときでも、予測に基づく行動決定が可能であり、学習が進むにつれて、徐々に、政策に基づく行動決定の割合が高くなっていく。

(第3の実施形態)

上述の実施形態では、予測に基づく行動決定部 16, 16 A の機能は予め与えられているものとしたが、行動生成機能を予め与えることが難しい場合は、目標状態に到達するための行動生成能力の獲得が必要になる。本実施形態では、このような行動生成能力を学習によって獲得させるものとする。

図 9 は図 6 の構成における予測に基づく行動決定部 16 A の本実施形態に係る内部構成を示す図である。図 9 において、31 は状態値 $s(t)$ を受け、この状態値 $s(t)$ が表す現在状態から、その前ステップにおける状態と行動を検出する行動付状態変化検出部、32 は行動付状態変化検出部 31 によって検出された現在状態 $s(t)$ 並びにその前ステップにおける状態 $s(t-1)$ および行動 $a(t-1)$ の組合せを状態変化として記憶する行動付状態変化記憶部、33 は行動付状態変化記憶部 32 から、現在状態から目標状態までの期間の状態変化の履歴を検索し、この検索結果を基にして行動を決定する行動計画部である。

図 10 のフローチャートを参照して、図 9 の予測に基づく行動決定部 16 A の動作を説明する。図 10 (a) は状態変化を蓄積する場合の動作、図 10 (b) は行動計画を行う場合の動作を示す。この 2 つの動作は、同時並行して実行できる。

状態変化を蓄積する場合、まず、環境 11 から状態 $s(t)$ を受ける (S31)。この現在の状態 $s(t)$ とそのときの行動 $a(t)$ はワーキングメモリに蓄えておく。そして、前ステップにおける状態 $s(t-1)$ と行動 $a(t-1)$ をワーキングメモリから取り出し (S32)、状態 $s(t)$ とともに行動付状態変化記憶部 32 に格納する (S33)。これは、状態 $s(t-1)$ のときに行動 $a(t-1)$ をとったら状態 $s(t)$ に変化したという、行動に起因する状態変化を表している。

また、行動計画を行う場合、まず、目標状態決定部 15 A から送られてきた目標状態を探索したい状態 $x_s(n)$ として設定する (S34)。そして、探索したい状態 $x_s(n)$ を行動付状態変化記憶部 32 から探索し (S35)、検出されたときは (S36 で YES)、探索したい状態 $x_s(n)$ と対になって記憶されている状態 $x_s(n-1)$ および行動 $x_a(n-1)$ をワーキングメモリに格納する (S37)。その後、ステップ S38 に進み、探索したい状態の 1 ステッ

ブ前の状態 $x_s(n-1)$ が現在の状態 $s(t-1)$ でなければ (S 3 8 で NO)、探索したい状態を更新し (S 3 9)、ステップ S 3 5 に戻る。同様の処理を繰り返し実行し、ステップ S 3 8 において、探索したい状態の 1 ステップ前の状態 $x_s(n-1)$ が現在の状態 $s(t-1)$ と一致したとき (YES)、それまでにワーキングメモリに格納されていた状態 x_s と行動 x_a の系列を行動計画として出力する (S 3 B)。

一方、ステップ S 3 6 において、探索したい状態 $x_s(n)$ が行動付状態変化記憶部 3 2 から検出できないときは (NO)、行動不能と判断して (S 3 A)、処理を終了する。なおこの場合は、予測に基づく行動決定部 1 6 A から、正しい行動が決定できないという信号が出力され、行動選択部 2 4 は、政策に基づく行動決定部 2 3 から出力された第 2 の行動候補を実際の行動として選択する。

このような動作によって、現在の行動のみならず、現在状態から目標状態に至るまでの行動計画が得られるので、行動計画が一旦完了した後は、その行動計画に従って行動候補を順に出力すればよい。これにより、処理量が格段に少なくなるので、特に長期にわたる予測誤差が少ない場合は、好ましい。もちろん、毎ステップごとに、目標状態までの行動計画を算出し直してもよい。この場合は、予測が完全でない場合であっても、行動が決定できる。

なお、本実施形態では、目標状態から現在状態に向かって後ろ向き探索を行うものとしたが、現在の状態 $s(t)$ と行動 $a(t)$ から $s(t+1)$ を算出する前向き探索を用いても、同様に行動計画を作成することができる。

(第 4 の実施形態)

本実施形態では、学習によって状態予測を実現するものとする。

図 1 1 は図 4 および図 6 の構成における環境予測部 1 3 の本実施形態に係る内部構成を示す図である。図 1 1 において、4 1 は状態値 $s(t)$ を受け、この状態値 $s(t)$ が表す現在状態から、その前ステップにおける状態を検出する状態変化検出部、4 2 は状態変化検出部 4 1 によって検出された現在状態 $s(t)$ およびその前ステップにおける状態 $s(t-1)$ を状態変化として記憶する状態変化記憶部、4 3 は状態変化記憶部 4 2 から、現在状態の後の状態を予測する状態予測部である。

図 12 のフローチャートを参照して、図 11 の環境予測部 13 の動作を説明する。図 12 (a) は状態変化を蓄積する場合の動作、図 12 (b) は状態予測を行う場合の動作を示す。この 2 つの動作は、同時並行して実行できる。

状態変化を蓄積する場合、まず、環境 10 から状態 $s(t)$ を受ける (S41)。この現在の状態 $s(t)$ はワーキングメモリに蓄えておく。そして、前ステップにおける状態 $s(t-1)$ の組合せをワーキングメモリから取り出し (S42)、状態 $s(t)$ とともに状態変化記憶部 42 に格納する。これは、状態 $s(t-1)$ の次に状態 $s(t)$ になったという状態変化を表している。

また、状態予測を行う場合、まず、環境 11 から取得した現在の状態 $s(t)$ を探索したい状態 $y_s(n)$ として設定する (S44)。そして、探索したい状態 $y_s(n)$ を状態変化記憶部 42 から探索し (S45)、検出されたときは (S46 で YES)、探索したい状態 $y_s(n)$ と対になって記憶されている 1 ステップ後の状態 $y_s(n+1)$ を状態変化記憶部 42 から取り出し、出力する (S47)。その後、ステップ S48 に進み、価値変換部 21 を経て、目標状態決定部 15、15A から評価信号による再予測依頼を受けたときは (YES)、探索したい状態を更新し (S49)、ステップ S45 に戻る。

一方、ステップ S46 において、探索したい状態 $y_s(n)$ が状態変化記憶部 42 から検出できないときは (NO)、予測先不明と判断して (S4A)、処理を終了する。なおこの場合は、環境予測部 13 から、正しい予測ができないという信号が出力され、行動選択部 24 は、政策に基づく行動決定部 23 から出力された第 2 の行動候補を実際の行動として選択する。

このような方法によって、環境予測部 13 の機能を予め作成しなくても、学習によって取得することができる。

なお、状態変化記憶部 42 の学習のために、ニューラル・ネットワーク等の関数近似手法によって次の状態を予測させることもできる。この場合、ニューラル・ネットワークが本来持つ汎化能力によって、経験したことがない状態 $s(t)$ に対しても、適切な 1 ステップ後の状態 $s(t+1)$ を出力できる可能性がある。

なお、上述の各実施形態では、主に、ピンポンのようなボールを打つ課題を例

に挙げて説明を行ったが、本発明の応用先は、このようなボールを打つ課題に限定されるものではなく、例えば、知能化住宅の制御、情報端末のソフトウェアエージェントの動作決定、ロボットの移動や行動決定など様々な用途が考えられる。

そして、本発明の特徴の1つは、予測した将来状態の中から、状態価値に基づいて、自律的に目標状態を決定することができる点にあり、状態価値が変わると、これに応じて目標状態も自動的に変化し得る。

例えば、知能化住宅においてエアコン等の機器の制御を行い、室内温度を快適にコントロールする場合、ボールの軌跡と同様に、外気温の影響を受けて室内温度がどのように変化するかを予測することは可能である。このとき、将来状態に対する状態価値がすでに学習されている場合（例えば、家の人が帰宅する時刻に近いほど、室内温度が快適になっている方が価値が高い）は、将来状態の中から目標となる時刻と室温の状態を決定することができる。そして、決定した目標状態（帰宅時刻と室温）に向けて、事前に機器の制御を行うことができる。気温の制御等、制御の効果がすぐには現れないような機器制御では、本発明のような予測に基づいた行動決定はきわめて有効である。

また、情報端末でスケジュール管理がなされている場合、ボールの軌跡と同様に、ユーザの将来の予定を参照することができる。このとき、情報提示やサービス提供などの行動（例えば、経済ニュースのダイジェストを提示するサービス）について、将来状態に対して状態価値がすでに学習されている場合（例えば、出張先への移動中が状態価値が高い）、目標状態を決定し、決定した目標状態（移動時刻）に向けて、事前の動作（ニュースの検索、ダウンロード、編集）を開始できる。スケジュールに応じて各将来状態の状態価値が変化すると、それに応じて目標状態も自動的に変化し、適切なタイミングで情報提示やサービス提供を行うことができる。

産業上の利用可能性

本発明では、将来の状態変化がより適切に考慮され、行動決定の精度が向上するので、例えば、知能化住宅の制御、情報端末のソフトウェアエージェントの動

作決定、家庭用ロボットの制御技術などに有効である。

請求の範囲

1. 所定の環境について状態を観察し、状態値を取得する状態観察部と、
前記環境の各状態に係る状態価値を記憶する状態価値記憶部と、
前記状態観察部によって取得された状態値に基づいて、前記環境の将来の状態
変化を予測する環境予測部と、

前記環境予測部によって予測された各将来状態の中から、前記状態価値記憶部
に格納された前記各将来状態の状態価値に基づいて、行動決定のために適した将
来の状態を、目標状態として決定する目標状態決定部と、

前記目標状態決定部によって決定された目標状態を基にして、自己の行動を決
定する第1の行動決定部とを備えた
ことを特徴とする予測型行動決定装置。

2. 請求項1において、

前記環境予測部は、前記環境の将来の、自己の行動に影響されない状態変化を
予測するものである
ことを特徴とする予測型行動決定装置。

3. 請求項1において、

前記目標状態決定部は、状態価値が極大となる将来状態を、目標状態として決
定する
ことを特徴とする予測型行動決定装置。

4. 請求項1において、

前記状態価値記憶部に記憶された状態価値を、学習によって更新する状態価値
更新部を備え、

前記目標状態決定部は、前記各将来状態の中で、状態価値が前記状態価値更新
部によってすでに更新されているものを、前記目標状態として決定する
ことを特徴とする予測型行動決定装置。

5. 請求項 1 において、

前記目標状態決定部は、前記状態価値記憶部から得た状態価値を、現在からのステップ数に応じて割り引いて、用いることを特徴とする予測型行動決定装置。

6. 請求項 1 において、

前記状態価値記憶部は、自己を含む状態に係る状態価値を記憶するものであり、

当該予測型行動決定装置は、

前記環境予測部によって予測された、自己を含まない将来状態について、前記状態価値記憶部に記憶された状態価値を基にしてその状態価値を求め、前記目標状態決定部に与える価値変換部を備えたものであることを特徴とする予測型行動決定装置。

7. 請求項 1 において、

所定の行動基準に基づいて、自己の行動を決定する第 2 の行動決定部と、

前記第 1 および第 2 の行動決定部によって決定された行動を第 1 および第 2 の行動候補として受け、これら第 1 および第 2 の行動候補のうちのいずれか一方を、実際の行動として選択する行動選択部とを備えたことを特徴とする予測型行動決定装置。

8. 請求項 7 において、

前記目標状態決定部は、目標状態を決定できたか否かを示す選択信号を前記行動選択部に与えるものであり、

前記行動選択部は、前記選択信号が、目標状態を決定できたことを示すときは、前記第 1 の行動候補を選択する一方、目標状態を決定できなかったことを示すときは、前記第 2 の行動候補を選択するものであることを特徴とする予測型行動決定装置。

9. 請求項1において、

前記第1の行動決定部は、

前記状態値を受け、この状態値が表す現在状態から、その前ステップにおける状態と行動を検出する行動付状態変化検出部と、

前記行動付状態変化検出部によって検出された、現在状態並びにその前ステップにおける状態および行動の組合せを、状態変化として記憶する行動付状態変化記憶部と、

前記行動付状態変化記憶部から、現在状態から目標状態までの期間の状態変化の履歴を検索し、この検索結果を基にして、行動を決定する行動計画部とを備えたものである

ことを特徴とする予測型行動決定装置。

10. 請求項9において、

前記行動計画部は、前記状態変化記憶部の検索の際に、目標状態から現在の状態に向かって後ろ向き探索を行うものである

ことを特徴とする予測型行動決定装置。

11. 請求項1において、

前記環境予測部は、

前記状態値を受け、この状態値が表す現在状態から、その前ステップにおける状態を検出する状態変化検出部と、

前記状態変化検出部によって検出された、現在状態およびその前ステップにおける状態の組合せを、状態変化として記憶する状態変化記憶部と、

前記状態変化記憶部から、現在状態の後の状態を予測する状態予測部とを備えたものである

ことを特徴とする予測型行動決定装置。

12. 予測型行動決定装置において、自己の行動を決定する方法であって、

所定の環境について状態を観察して、状態値を取得する第1のステップと、
取得した状態値に基づいて、前記環境の将来の状態変化を予測する第2のステップと、

予測した各将来状態の中から、前記各将来状態の状態価値を参照しつつ、行動決定のために適した将来状態を、目標状態として決定する第3のステップと、

決定した目標状態を基にして、自己の行動を決定する第4のステップとを備えた
ことを特徴とする行動決定方法。

13. 請求項12において、

予測する状態変化は、前記環境の将来の、自己の行動に影響されない状態変化である

ことを特徴とする行動決定方法。

14. 請求項12において、

前記第3のステップにおいて、

前記各将来状態の中で、状態価値が極大となるものを、前記目標状態として決定する

ことを特徴とする行動決定方法。

15. 請求項12において、

前記予測型行動決定装置は、前記環境の各状態に係る状態価値を、学習によって更新するものであり、

前記第3のステップにおいて、

前記各将来状態の中で、状態価値がすでに更新されているものを、前記目標状態として決定する

ことを特徴とする行動決定方法。

補正書の請求の範囲

[2004年6月03日(03.06.04)国際事務局受理：出願当初の請求の範囲2及び13は取り下げられた；出願当初の請求の範囲1, 5, 11及び12は補正された；他の請求の範囲は変更なし。(5頁)]

1. (補正後) 所定の環境について状態を観察し、状態値を取得する状態観察部と、

前記環境の各状態に係る状態価値を記憶する状態価値記憶部と、

前記状態観察部によって取得された状態値に基づいて、前記環境の将来の状態変化を予測する環境予測部と、

前記環境予測部によって予測された各将来状態の中から、前記状態価値記憶部に格納された前記各将来状態の状態価値に基づいて、行動決定のために適した将来の状態を、目標状態として決定する目標状態決定部と、

前記目標状態決定部によって決定された目標状態を基にして、自己の行動を決定する第1の行動決定部とを備え、

前記環境予測部は、前記環境の将来の、自己の行動に影響されない状態変化を予測するものである

ことを特徴とする予測型行動決定装置。

2. (削除)

3. 請求項1において、

前記目標状態決定部は、状態価値が極大となる将来状態を、目標状態として決定する

ことを特徴とする予測型行動決定装置。

4. 請求項1において、

前記状態価値記憶部に記憶された状態価値を、学習によって更新する状態価値更新部を備え、

前記目標状態決定部は、前記各将来状態の中で、状態価値が前記状態価値更新部によってすでに更新されているものを、前記目標状態として決定する

ことを特徴とする予測型行動決定装置。

5. (補正後) 所定の環境について状態を観察し、状態値を取得する状態観察部と、

前記環境の各状態に係る状態価値を記憶する状態価値記憶部と、

前記状態観察部によって取得された状態値に基づいて、前記環境の将来の状態変化を予測する環境予測部と、

前記環境予測部によって予測された各将来状態の中から、前記状態価値記憶部に格納された前記各将来状態の状態価値に基づいて、行動決定のために適した将来の状態を、目標状態として決定する目標状態決定部と、

前記目標状態決定部によって決定された目標状態を基にして、自己の行動を決定する第1の行動決定部とを備え、

前記目標状態決定部は、前記状態価値記憶部から得た状態価値を、現在からのステップ数に応じて割り引いて、用いることを特徴とする予測型行動決定装置。

6. 請求項1において、

前記状態価値記憶部は、自己を含む状態に係る状態価値を記憶するものであり、

当該予測型行動決定装置は、

前記環境予測部によって予測された、自己を含まない将来状態について、前記状態価値記憶部に記憶された状態価値を基にしてその状態価値を求め、前記目標状態決定部に与える価値変換部を備えたものであることを特徴とする予測型行動決定装置。

7. 請求項1において、

所定の行動基準に基づいて、自己の行動を決定する第2の行動決定部と、

前記第1および第2の行動決定部によって決定された行動を第1および第2の行動候補として受け、これら第1および第2の行動候補のうちのいずれか一方を、実際の行動として選択する行動選択部とを備えた

ことを特徴とする予測型行動決定装置。

8. 請求項7において、

前記目標状態決定部は、目標状態を決定できたか否かを示す選択信号を前記行動選択部に与えるものであり、

前記行動選択部は、前記選択信号が、目標状態を決定できたことを示すときは、前記第1の行動候補を選択する一方、目標状態を決定できなかったことを示すときは、前記第2の行動候補を選択するものであることを特徴とする予測型行動決定装置。

9. 請求項1において、

前記第1の行動決定部は、

前記状態値を受け、この状態値が表す現在状態から、その前ステップにおける状態と行動を検出する行動付状態変化検出部と、

前記行動付状態変化検出部によって検出された、現在状態並びにその前ステップにおける状態および行動の組合せを、状態変化として記憶する行動付状態変化記憶部と、

前記行動付状態変化記憶部から、現在状態から目標状態までの期間の状態変化の履歴を検索し、この検索結果を基にして、行動を決定する行動計画部とを備えたものである

ことを特徴とする予測型行動決定装置。

10. 請求項9において、

前記行動計画部は、前記状態変化記憶部の検索の際に、目標状態から現在の状態に向かって後ろ向き探索を行うものである

ことを特徴とする予測型行動決定装置。

11. (補正後) 所定の環境について状態を観察し、状態値を取得する状態観察部と、

前記環境の各状態に係る状態価値を記憶する状態価値記憶部と、

前記状態観察部によって取得された状態値に基づいて、前記環境の将来の状態変化を予測する環境予測部と、

前記環境予測部によって予測された各将来状態の中から、前記状態価値記憶部に格納された前記各将来状態の状態価値に基づいて、行動決定のために適した将来の状態を、目標状態として決定する目標状態決定部と、

前記目標状態決定部によって決定された目標状態を基にして、自己の行動を決定する第1の行動決定部とを備え、

前記環境予測部は、

前記状態値を受け、この状態値が表す現在状態から、その前ステップにおける状態を検出する状態変化検出部と、

前記状態変化検出部によって検出された、現在状態およびその前ステップにおける状態の組合せを、状態変化として記憶する状態変化記憶部と、

前記状態変化記憶部から、現在状態の後の状態を予測する状態予測部とを備えたものである

ことを特徴とする予測型行動決定装置。

12. (補正後) 予測型行動決定装置において、自己の行動を決定する方法であって、

所定の環境について状態を観察して、状態値を取得する第1のステップと、

取得した状態値に基づいて、前記環境の将来の状態変化を予測する第2のステップと、

予測した各将来状態の中から、前記各将来状態の状態価値を参照しつつ、行動決定のために適した将来状態を、目標状態として決定する第3のステップと、

決定した目標状態を基にして、自己の行動を決定する第4のステップとを備え、

予測する状態変化は、前記環境の将来の、自己の行動に影響されない状態変化である

ことを特徴とする行動決定方法。

1 3. (削除)

1 4. 請求項 1 2 において、

前記第 3 のステップにおいて、

前記各将来状態の中で、状態価値が極大となるものを、前記目標状態として決定する

ことを特徴とする行動決定方法。

1 5. 請求項 1 2 において、

前記予測型行動決定装置は、前記環境の各状態に係る状態価値を、学習によって更新するものであり、

前記第 3 のステップにおいて、

前記各将来状態の中で、状態価値がすでに更新されているものを、前記目標状態として決定する

ことを特徴とする行動決定方法。

1/14

FIG. 1

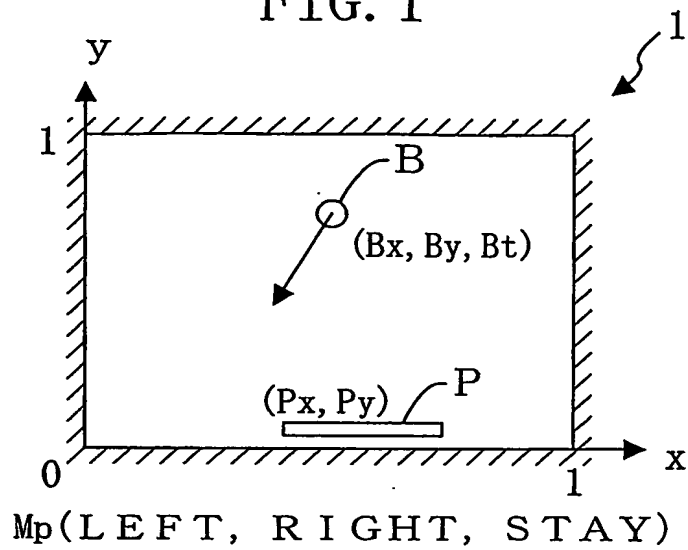
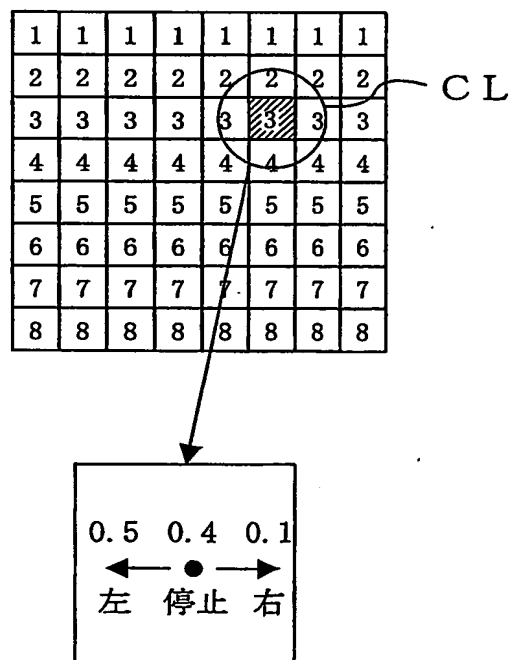


FIG. 2



2/14

FIG. 3

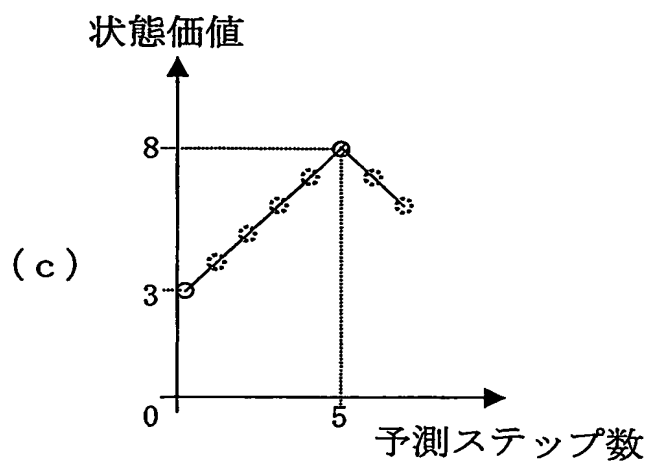
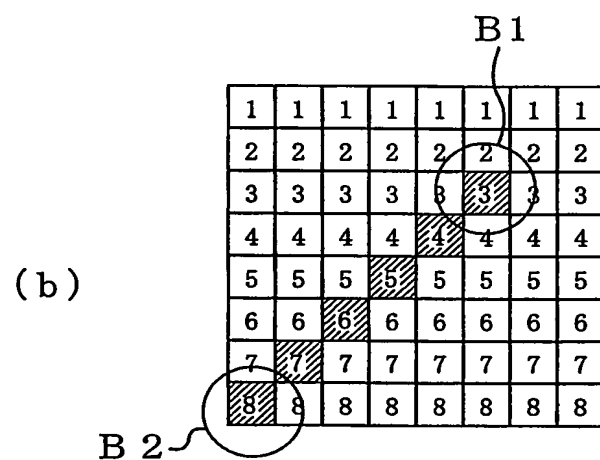
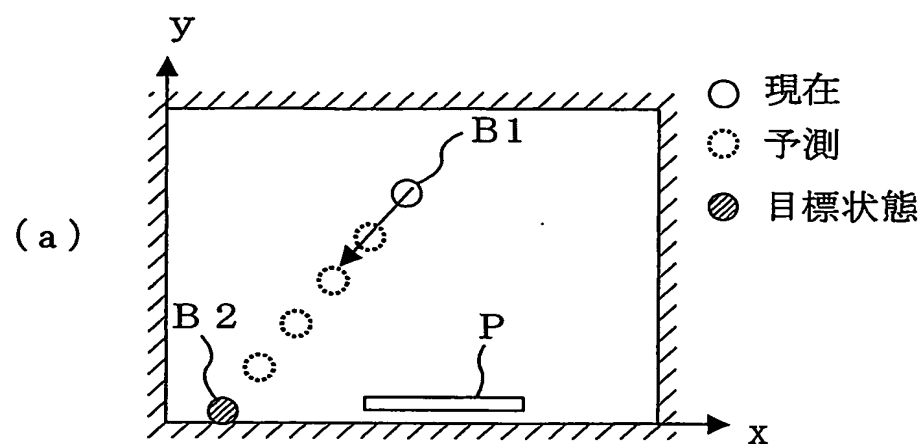
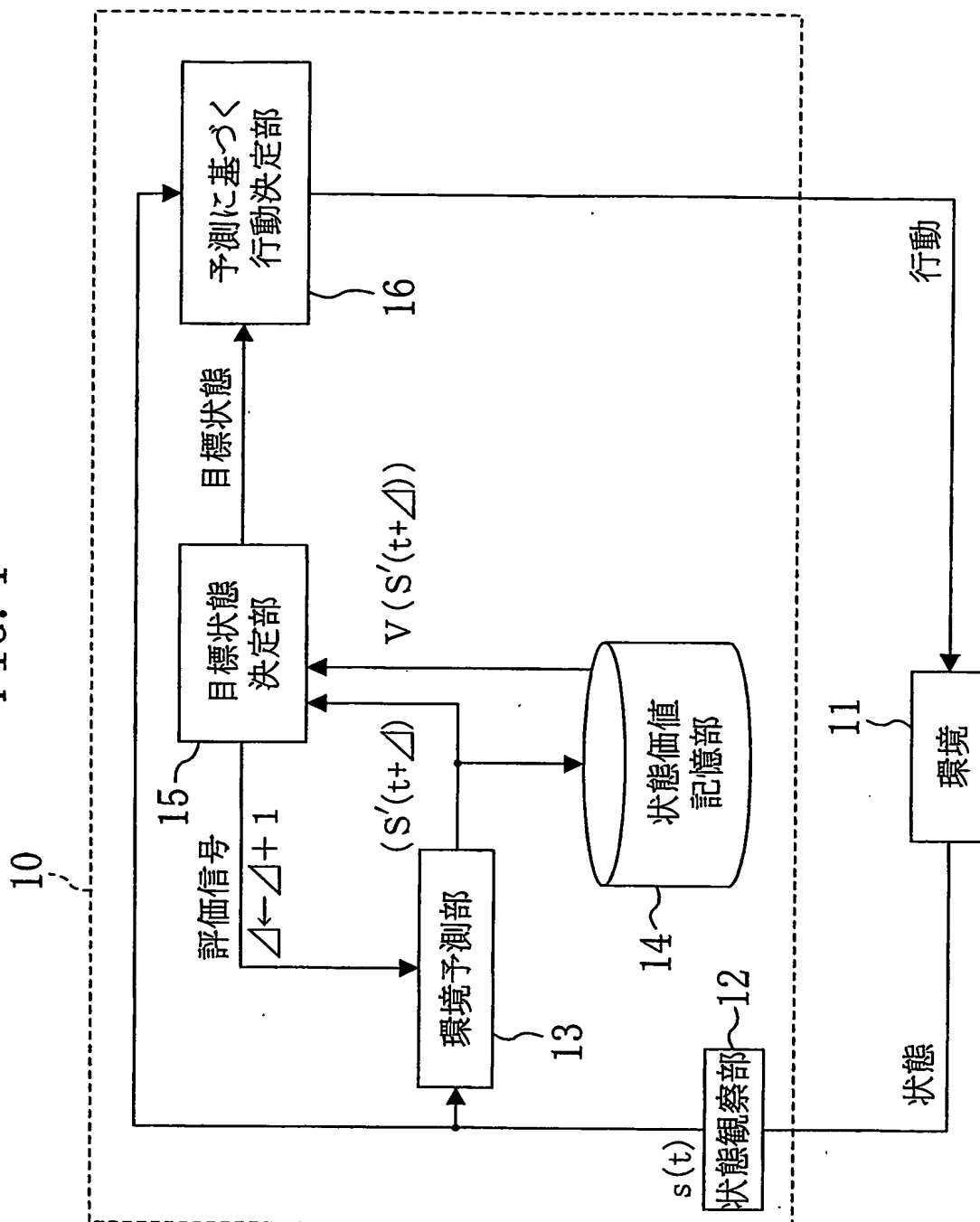


FIG. 4



4/14

FIG. 5

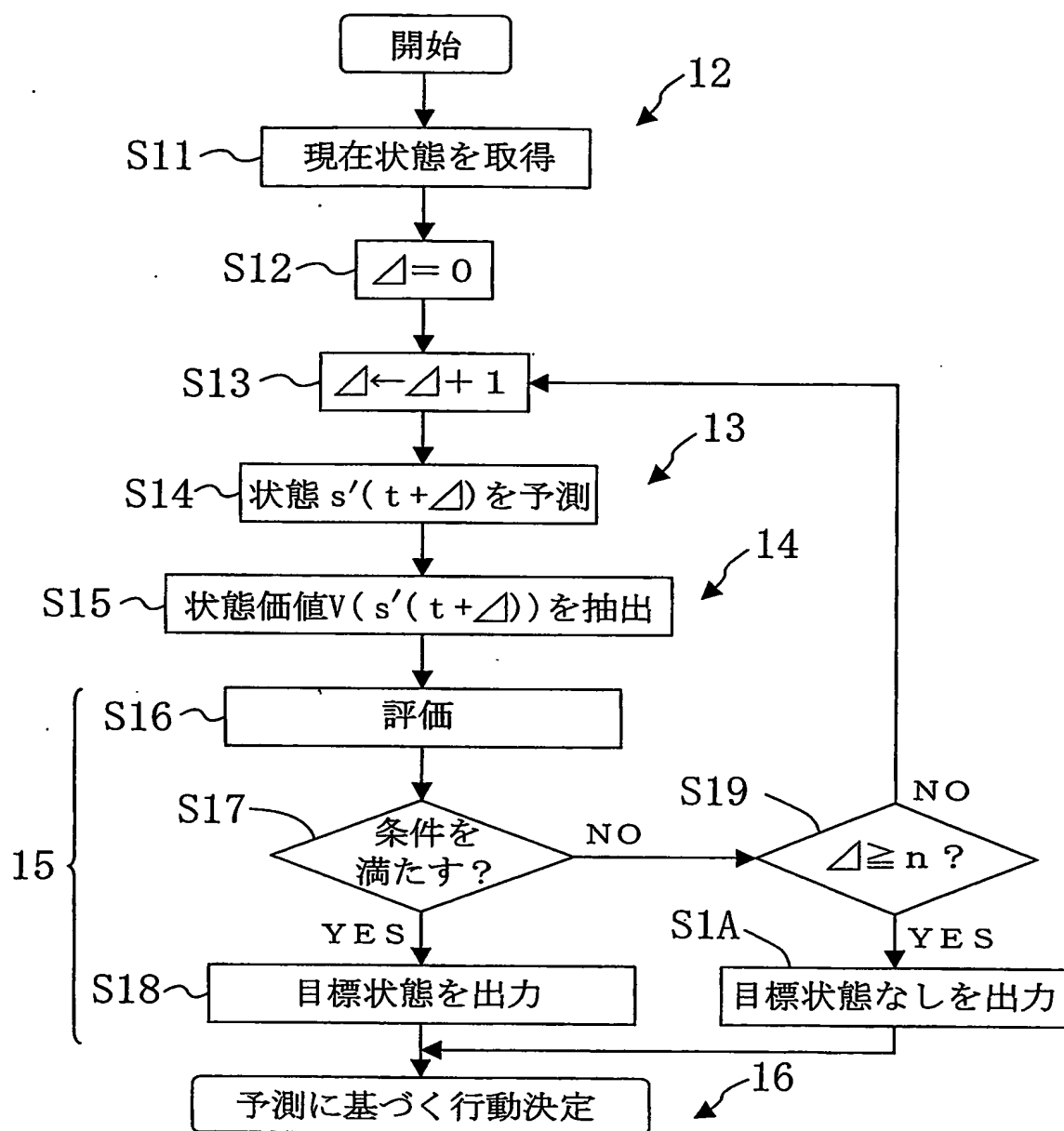
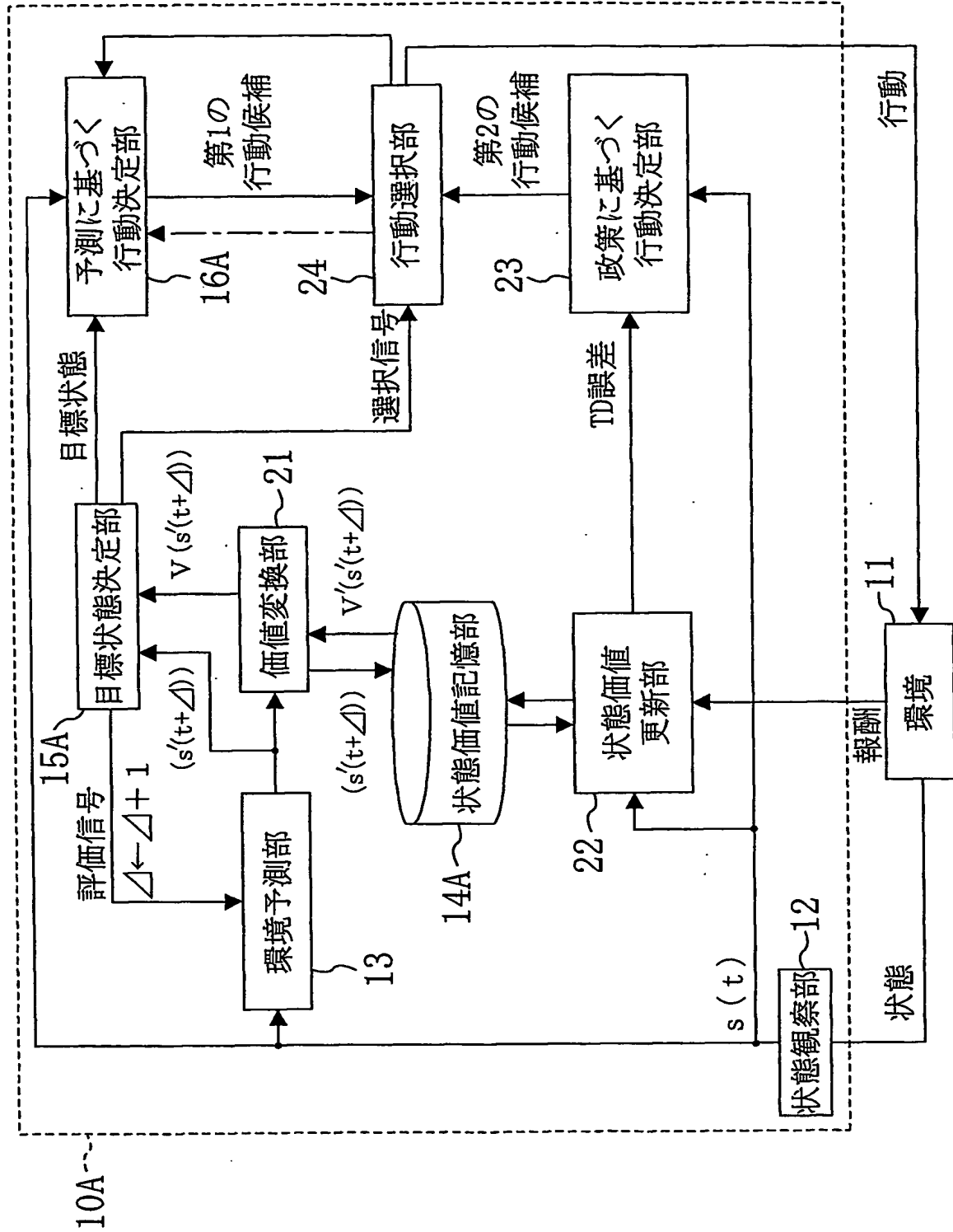
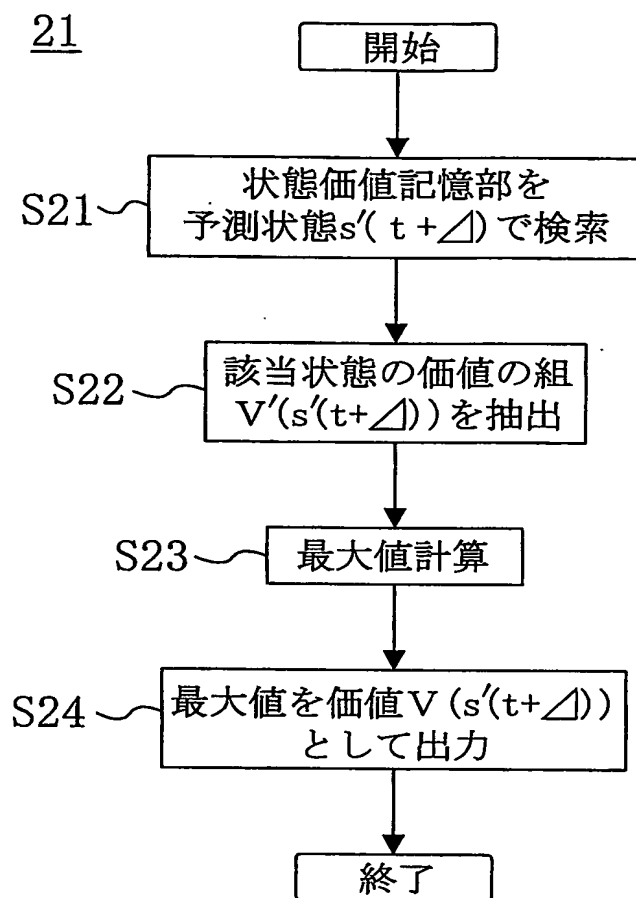


FIG. 6



6/14

FIG. 7



7/14

FIG. 8

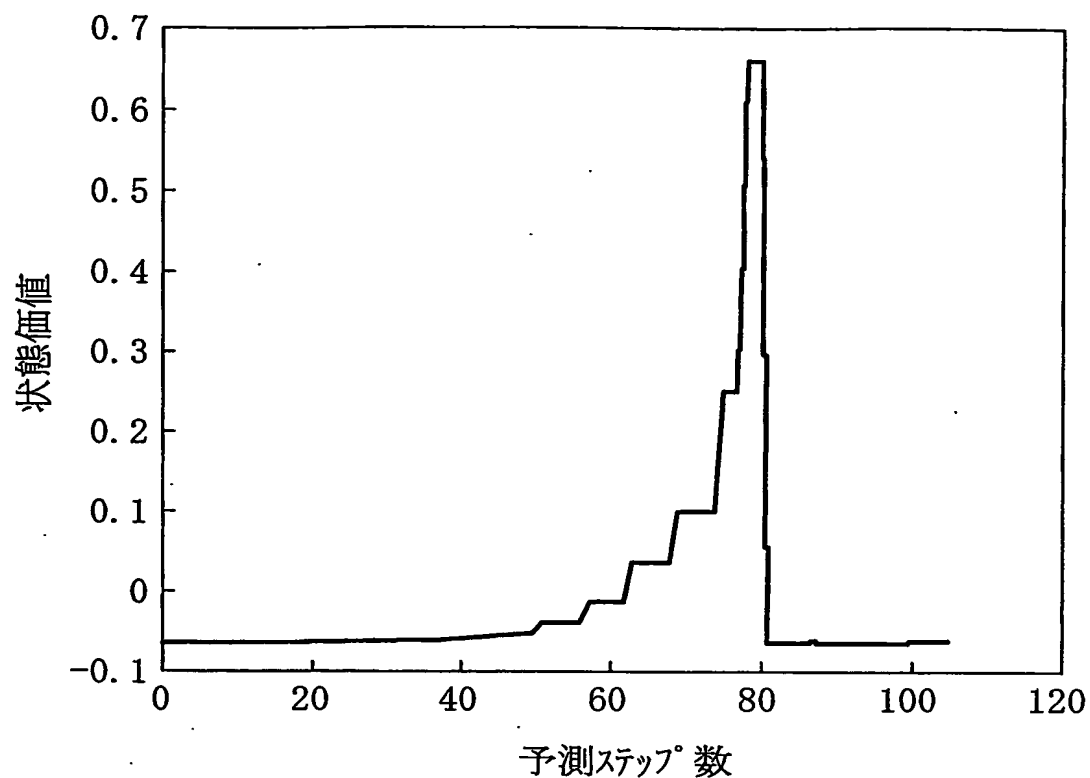


FIG. 9

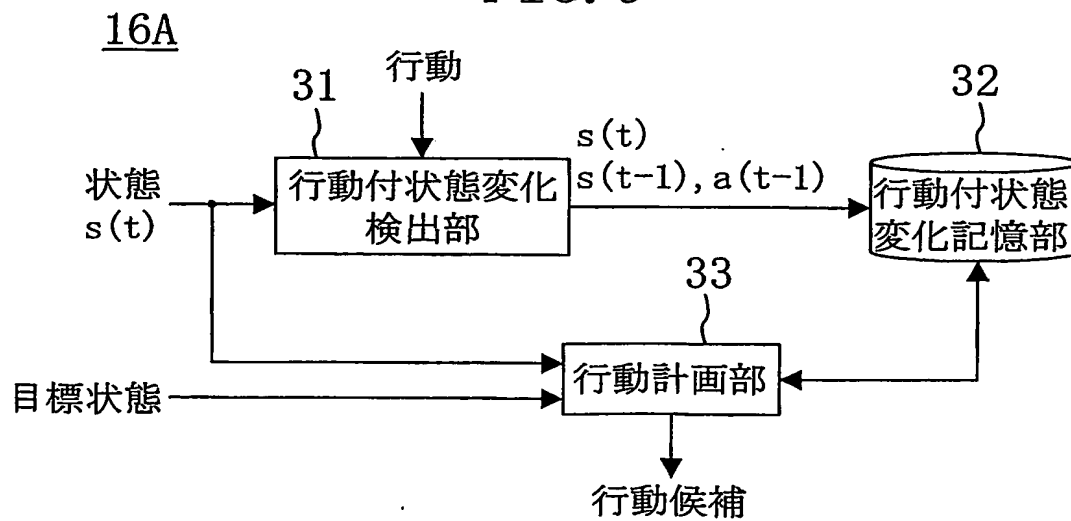
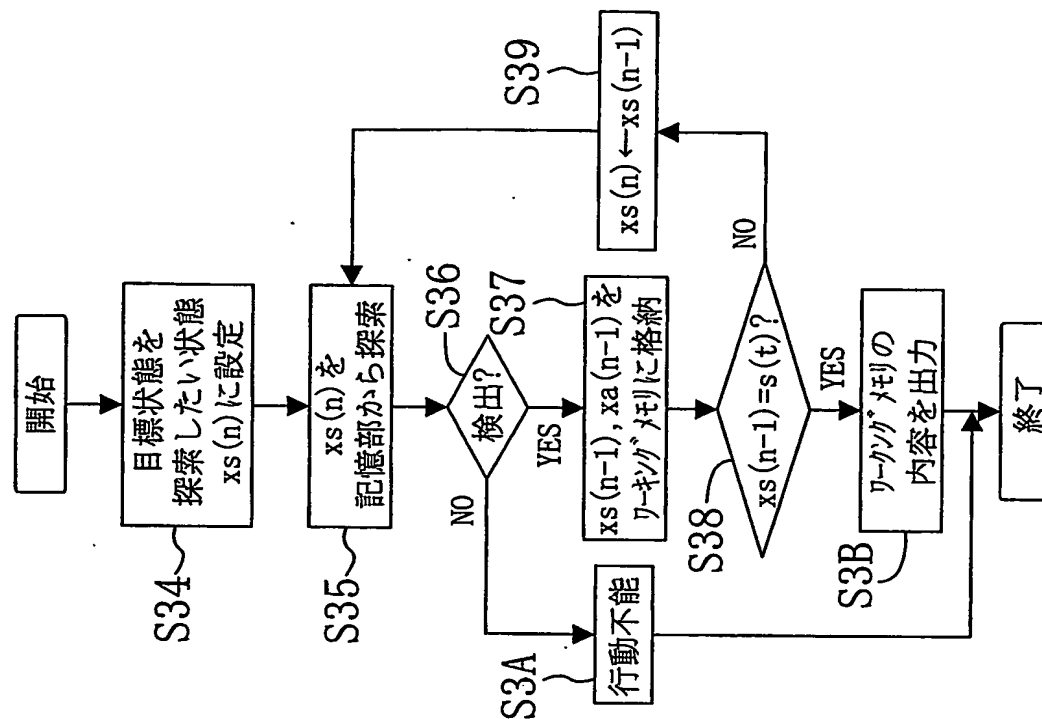


FIG. 10

(b) 行動計画



(a) 状態変化の蓄積

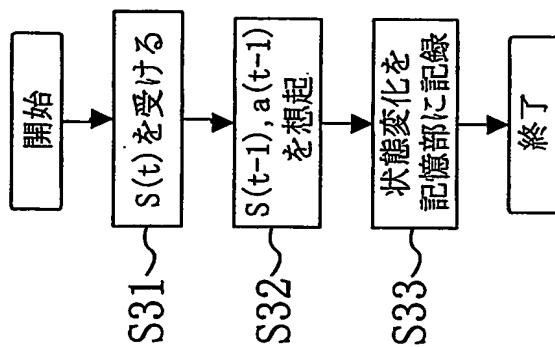


FIG. 11

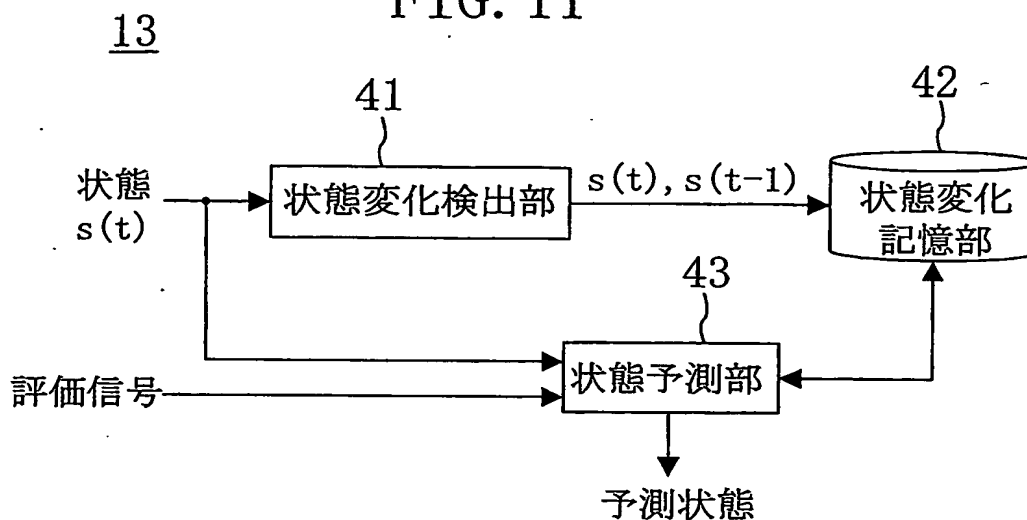
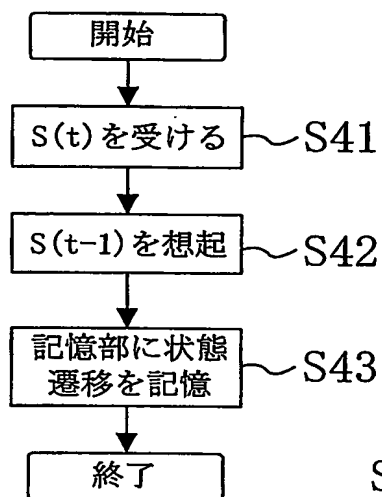


FIG. 12

(a) 状態遷移の蓄積



(b) 状態予測

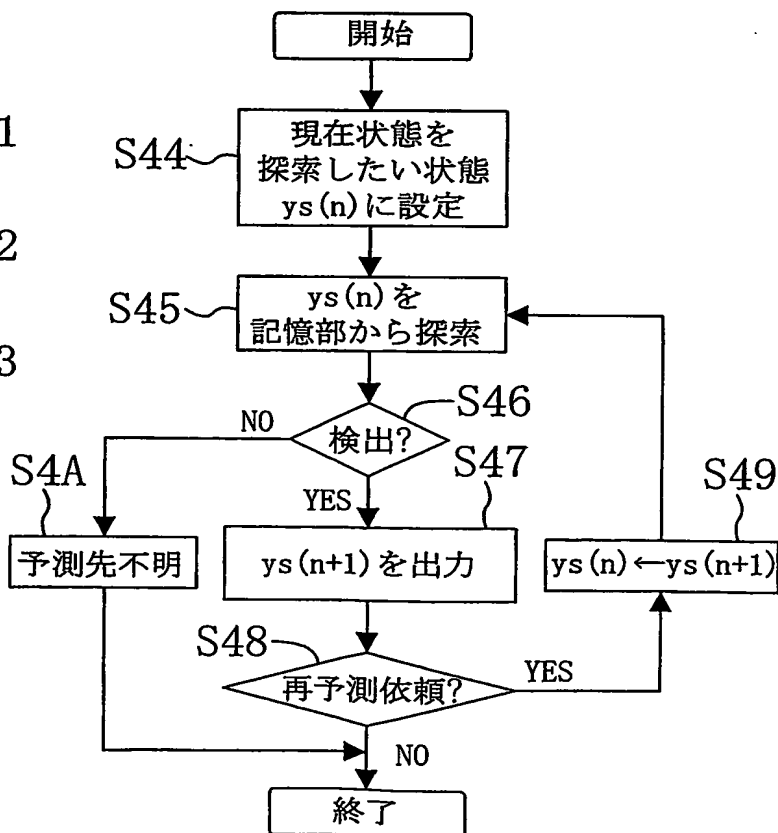


FIG. 13

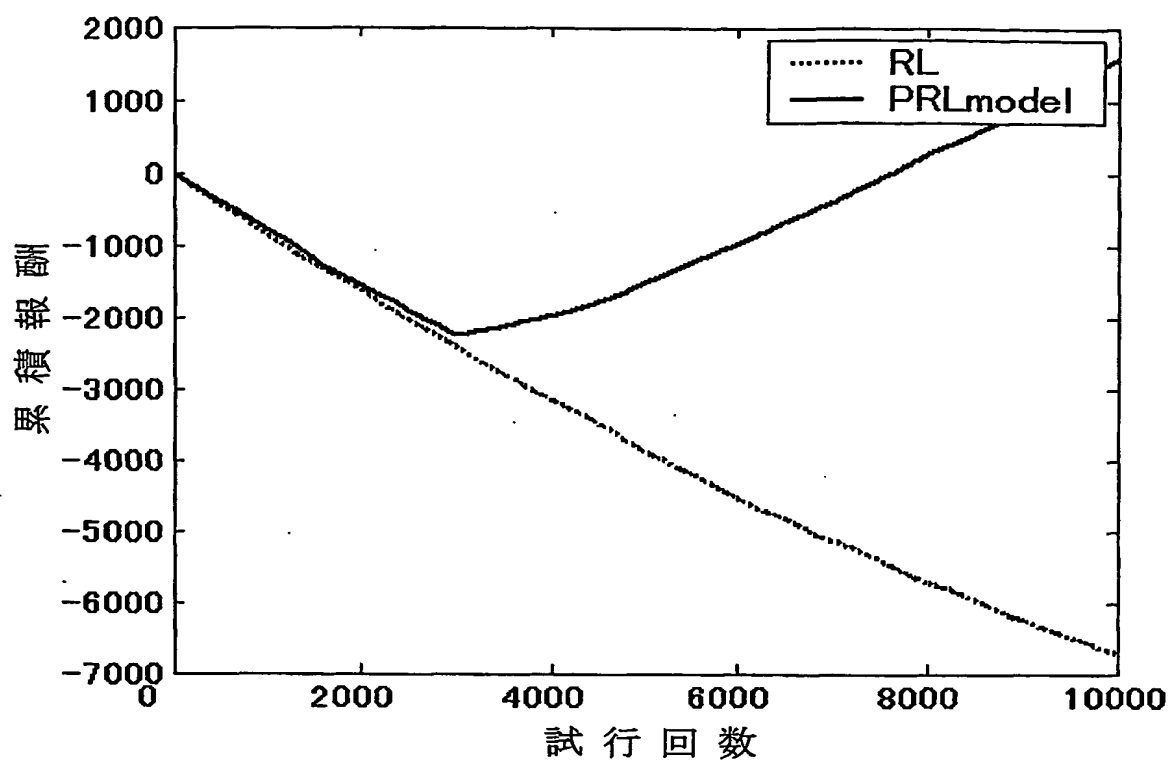
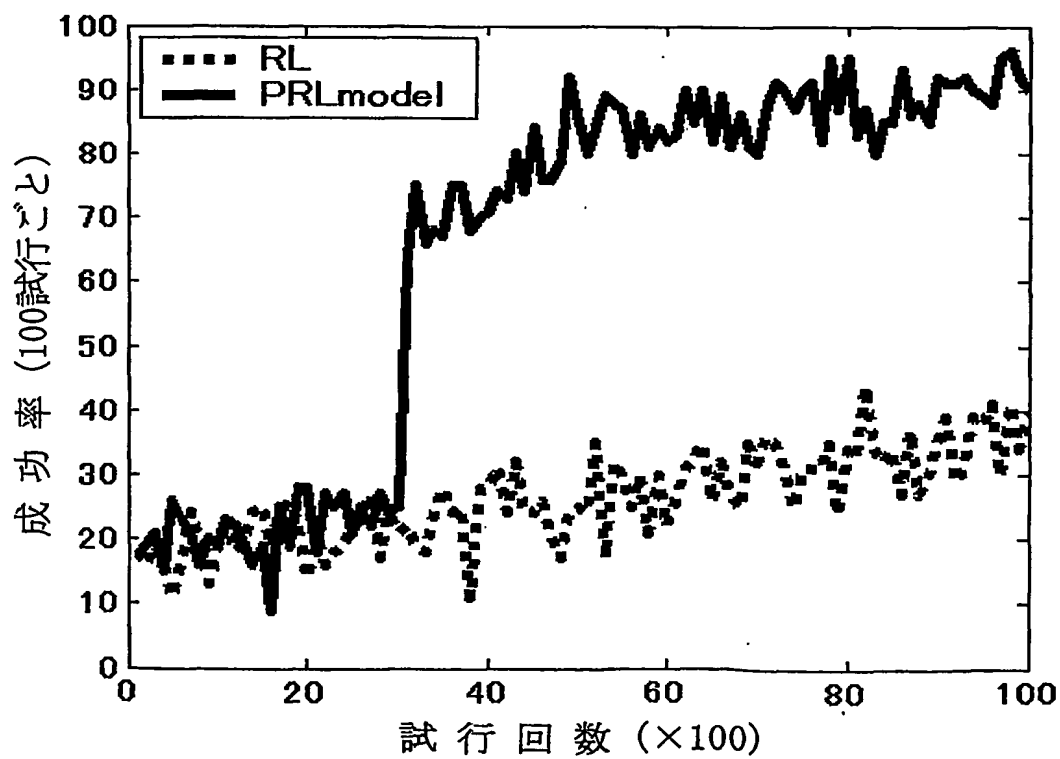
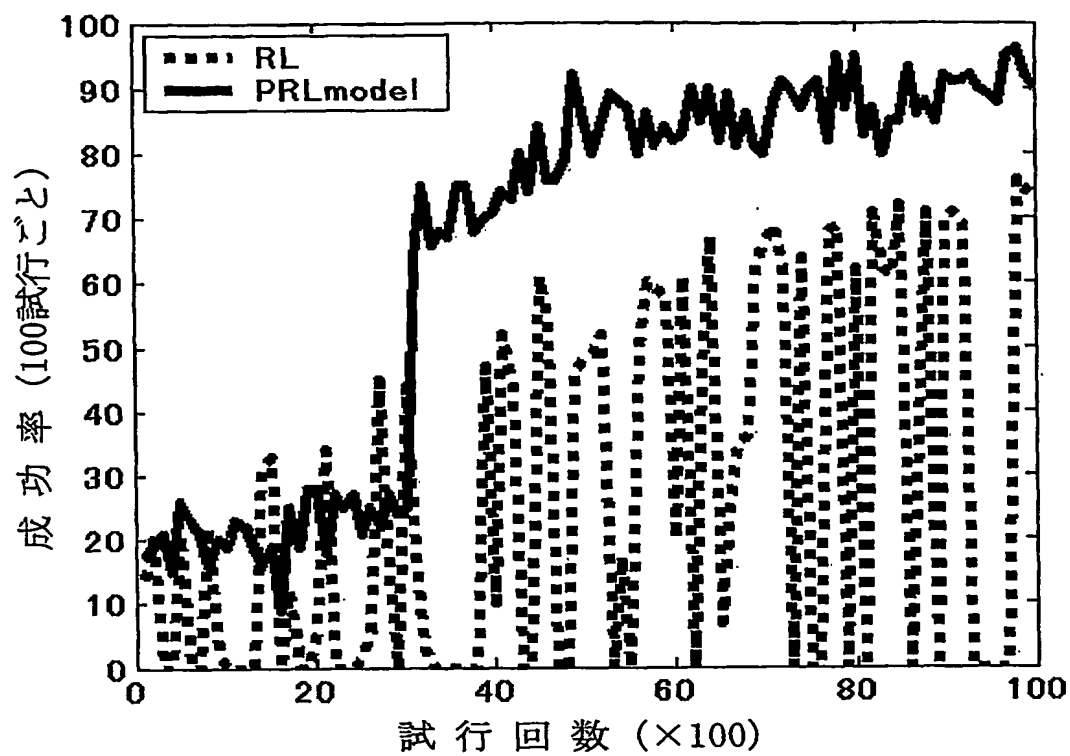


FIG. 14



11/14

FIG. 15



13/14

FIG. 17

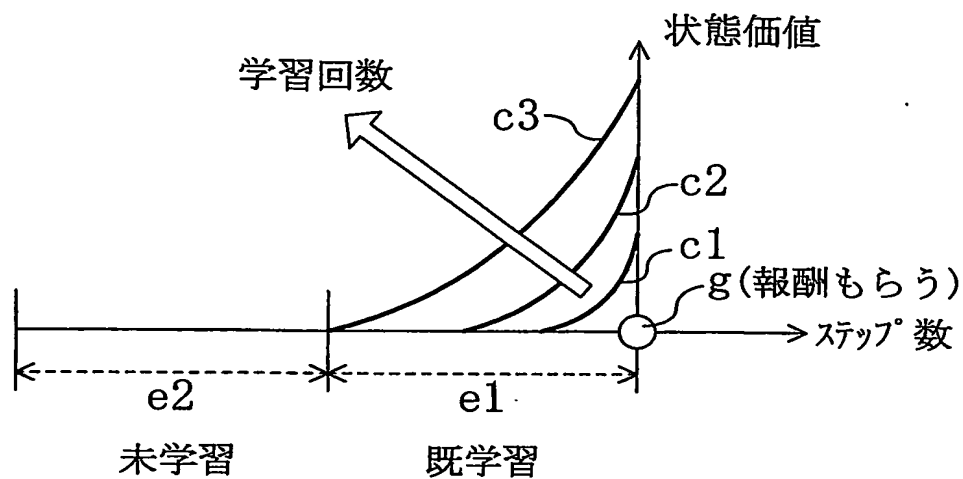
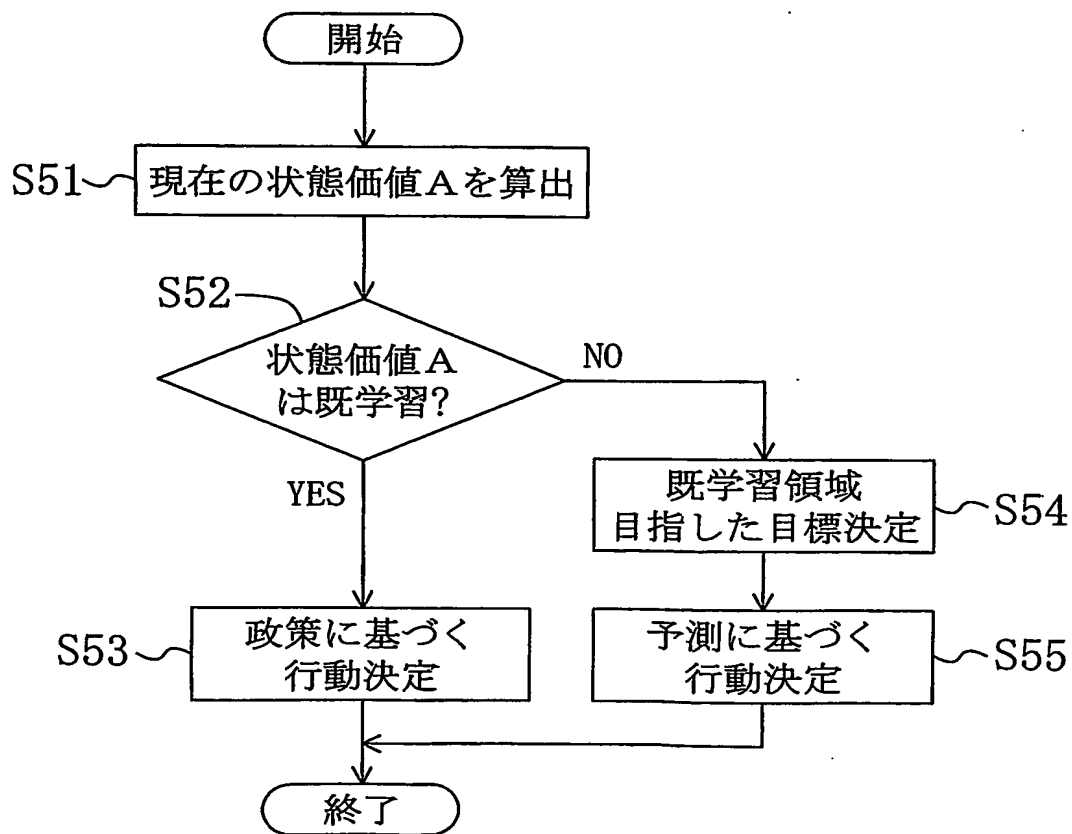
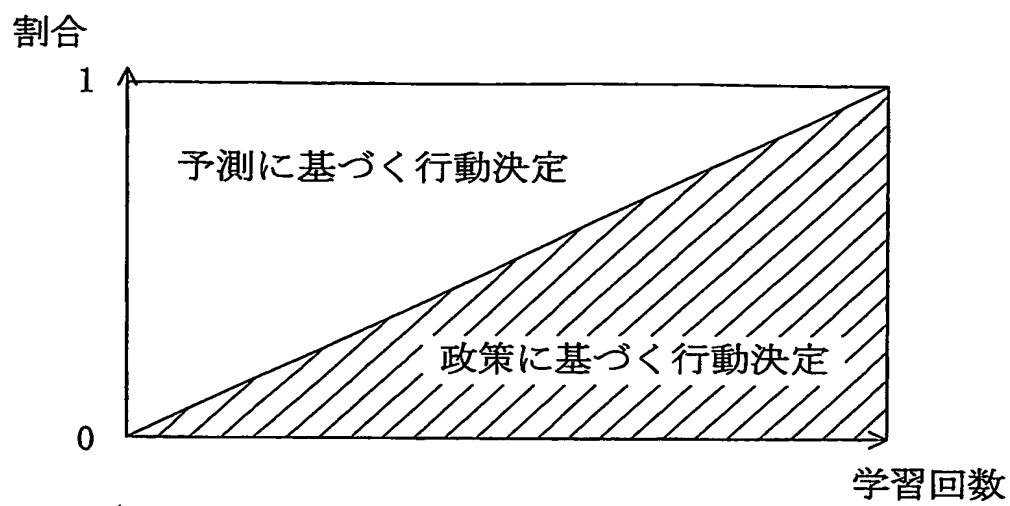


FIG. 18



14/14

FIG. 19



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2004/000971

A. CLASSIFICATION OF SUBJECT MATTER

Int.Cl.⁷ G06N5/00

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

Int.Cl.⁷ G06N5/00

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho	1922-1996	Toroku Jitsuyo Shinan Koho	1994-2004
Kokai Jitsuyo Shinan Koho	1971-2004	Jitsuyo Shinan Toroku Koho	1996-2004

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

JOIS, DAITO MASARU+OMORI TAKASHI+MORIKAWA KOJI+OKA NATSUKI
(in Japanese)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JP 9-297690 A (Nippon Telegraph And Telephone Corp.), 18 November, 1997 (18.11.97), Column 2, line 35 to column 3, line 39; Fig. 9 (Family: none)	1, 3, 12, 14
Y	same as the above	4, 7-10, 15
A	same as the above	2, 5, 6, 11, 13
Y	Taku YOSHIOKA et al., "Kyoka Gakushu ni yoru Othello no Senryaku Kakutoku", Dai 32. Kai Jinko Chino Kisoron Kenkyukai Shiryo, The Japanese Society for Artificial Intelligence, Japan, 26 March, 1998 (26.03.98), pages 115 to 120; page 116, left column, lines 22 to 35	4, 7, 8, 15

☒ Further documents are listed in the continuation of Box C.☐ See patent family annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search

24 March, 2004 (24.03.04)

Date of mailing of the international search report

06 April, 2004 (06.04.04)

Name and mailing address of the ISA/
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2004/000971

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP 9-245015 A (Toshiba Corp.), 19 September, 1997 (19.09.97), Column 20, lines 19 to 26; Fig. 6 (Family: none)	9, 10
Y	Takeshi ITO, "Shogi no Kansosen ni Mirareru Tansakuno Henka Katei", Information Processing Society of Japan Kenkyu Hokoku, Information Processing Society of Japan, Vol.2001, No.28, 16 March, 2001 (16.03.01), pages 47 to 54; page 48, left column, lines 31 to 40	10
A	Yasuharu KOIKE, Kenji Dotani, "Multi Step Jotai Yosoku o Mochiita Kyoka Gakushu-Chumokuten Idoten o Koryo shita Jidosha Unten Model", The Institute of Electronics, Information and Communication Engineers Gijutsu Kenkyu Hokoku, The Institute of Electronics, Information and Communication Engineers, Vol.98, No.673, 18 March, 1999 (18.03. 99), pages 241 to 247	1-15
A	Koichiro TAKITA et al., "Pulse Neural Network ni Okeru Network Kakuchogata Kyoka Gakushu Algorithm", The Institute of Electronics, Information and Communication Engineers Gijutsu Kenkyu Hokoku, The Institute of Electronics, Information and Communication Engineers, Vol.99, No.684, 13 March, 2000 (13.03.00), pages 181 to 188, full text; all drawings	1-15
A	Shingo OCHIAI et al., "Kyoka Gakushu o Mochiita Kyugi Dosa no Yosoku to Seigyo", The Institute of Electronics, Information and Communication Engineers Gijutsu Kenkyu Hokoku, The Institute of Electronics, Information and Communication Engineers, Vol.97, No.69, 23 May, 1997 (23.05.97), pages 69 to 75; full text; all drawings	1-15

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int. Cl⁷ G06N5/00

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int. Cl⁷ G06N5/00

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報 1922-1996年

日本国公開実用新案公報 1971-2004年

日本国登録実用新案公報 1994-2004年

日本国実用新案登録公報 1996-2004年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

JOIS, 大東優+大森隆司+森川幸治+岡夏樹

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
	JP 9-297690 A (日本電信電話株式会社) 1997. 11. 18	
X	第2欄第35行-第3欄第39行, 図9 (ファミリーなし)	1, 3, 12, 14
Y	同上	4, 7-10, 15,

☒ C欄の続きにも文献が列挙されている。☐ パテントファミリーに関する別紙を参照。

* 引用文献のカテゴリー

「A」 特に関連のある文献ではなく、一般的技術水準を示すもの

「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの

「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)

「O」 口頭による開示、使用、展示等に言及する文献

「P」 国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの

「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの

「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの

「&」 同一パテントファミリー文献

国際調査を完了した日

24. 03. 2004

国際調査報告の発送日

06. 4. 2004

国際調査機関の名称及びあて先

日本国特許庁 (ISA/JP)

郵便番号100-8915

東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

漆原 孝治

5B

3145

電話番号 03-3581-1101 内線 3546

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
A	同上	2, 5, 6, 11, 13
	吉岡琢, 外2名著, “強化学習によるオセロの戦略獲得”, 第32 回人工知能基礎論研究会資料, 社団法人人工知能学会, 日本, 1998. 03. 26, 第115-120頁	
Y	第116頁左欄22-35行	4, 7, 8, 15
	JP 9-245015 A (株式会社東芝) 1997. 09. 19	
Y	第20欄第19-26行, 図6 (ファミリーなし)	9, 10
	伊藤毅志 著, “将棋の感想戦に見られる探索の変化過程”, 情報 処理学会研究報告, 社団法人情報処理学会, 第2001巻, 第28 号, 2001. 03. 16, 第47-54頁	
Y	第48頁左欄第31-40行	10
	小池康晴・銅谷賢治共著, “マルチステップ状態予測を用いた強化学習-注目点移動点を考慮した自動車運転モデル”, 電子情報通信 学会技術研究報告, 社団法人電子情報通信学会, 第98巻, 第67 3号, 1999. 03. 18, 第241-247頁	
A	全文, 全図	1-15
	瀧田航一朗, 外2名著, “パルスニューラルネットワークにおける ネットワーク拡張型強化学習アルゴリズム”, 電子情報通信学会技 術研究報告, 社団法人電子情報通信学会, 第99巻, 第684号, 2000. 03. 13, 第181-188頁	
A	全文, 全図	1-15
	落合真吾, 外2名著, “強化学習を用いた球技動作の予測と制御 ”, 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 第97巻, 第69号, 1997. 05. 23, 第69-75頁	
A	全文, 全図	1-15